
Making Batch Normalization Great in Federated Deep Learning

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Batch Normalization (BN) is commonly used in modern deep foundation models
2 to improve stability and speed up convergence in centralized training. In federated
3 learning (FL) with non-IID decentralized data, previous works observed that training
4 with BN could hinder performance due to the mismatch of the BN statistics
5 between training and testing. Group Normalization (GN) is thus more often used
6 in FL as an alternative to BN. In this paper, we identify a more fundamental issue
7 of BN in FL that makes BN inferior even with high-frequency communication be-
8 tween clients and servers. We then propose a frustratingly simple treatment, which
9 significantly improves BN and makes it outperform GN across a wide range of FL
10 settings. Along with this study, we also reveal an unreasonable behavior of BN in
11 FL. We find it quite robust in the low-frequency communication regime where FL
12 is commonly believed to degrade drastically. We hope that our study could serve
13 as a valuable reference for future practical usage and theoretical analysis in FL.

14 1 Introduction

15 Foundation models [1] are normally very deep neural networks (DNNs) trained via stochastic gradient
16 descent (SGD). FEDAVG [2] is arguably the most widely used training algorithm in an FL setting.
17 FEDAVG iterates between two steps: parallel local SGD at the clients, and global model aggregation
18 at the server. In the extreme case where global aggregation takes place after every local SGD step,
19 FEDAVG is very much equivalent to centralized SGD for training simple DNN models like multi-layer
20 perception [3, 4, 5, 6, 7]. Of course, due to the communication costs in practice, it is unlikely for
21 clients to communicate at such a high frequency. Many existing works have thus focused on how
22 to train DNNs at a lower communication frequency (e.g., once after local SGD for a few epochs),
23 especially under the challenging condition where the data across clients are non-IID [8, 9, 10, 11].

24 In this paper, we specifically focus on DNN models that contain Batch Normalization (BN) layers [12].
25 In centralized training, especially for deep feed-forward models like ResNet [13], BN has been widely
26 used to improve the stability of training and speed up convergence. In the literature on FL, however,
27 many of the previous experiments have focused on shallow ConvNets (CNN) without BN; only a
28 few works have particularly studied the usage of BN in FL [14, 15]. In [14], the authors pointed
29 out the mismatch between the feature statistics (i.e., the means and variances in BN) estimated on
30 non-IID local data (during training) and global data (during testing), and argued that this cannot
31 be addressed by using larger mini-batch sizes or other sampling strategies. Hsieh et al. [14] thus
32 proposed to replace BN with Group Normalization (GN) [16] and showed its superior performance
33 in some extreme non-IID settings. Such a solution has since been followed by a long non-exhaustive
34 list of later works [17, 18, 19, 20, 21, 22, 23, 24].

35 With that being said, replacing BN with GN in FL is more like an ad hoc solution rather than a
36 cure-all. First, in centralized training, BN typically outperforms GN empirically. Replacing BN with

37 GN in FL thus seems like a compromise. Second, several recent works [25, 26, 27, 11] have reported
 38 that BN is still better than GN in their specific FL settings. Third, changing the normalization
 39 layer may create a barrier between the communities of centralized learning and FL. To illustrate, in
 40 centralized training, many publicly available pre-trained checkpoints [28, 29] are based on popular
 41 CNN architectures even recent transformers [30] with BN; most understanding [31, 32, 33], empirical
 42 studies [34], and theoretical analysis [35] about normalization in DNNs are built upon BN rather
 43 than GN. These prior results may become hard to be referred to in the FL community.

44 Last but not least, after a careful study of recent works that
 45 reported poor performance of BN [36, 37, 38], we found that
 46 the huge gap between centralized learning and FL cannot be
 47 closed even if clients communicate right after *every* local SGD
 48 step. Such a finding sharply contradicts what is observed on
 49 DNNs without BN. In other words, the issue with applying BN
 50 in FL seems to be more fundamental than previously believed.

51 **Contributions.** Building upon these aspects, we strive to an-
 52 swer the following questions towards a more holistic under-
 53 standing of BN in FL, especially under non-IID settings.

- 54 1. Why does BN degrade so drastically in FL compared to
 55 centralized training or other normalizers? (section 3)
- 56 2. Is there a way to properly use BN in FL to bridge the per-
 57 formance gap w.r.t. centralized training? (section 4)
- 58 3. Is there a comfort zone and danger zone for BN (and other
 59 normalization methods) in FL? (Appendix B)

60 To begin with, we investigate several different perspectives to understand the issue of BN in FL,
 61 including BN statistic dynamics, the training/test mismatch of statistics, and the gradient w.r.t. the
 62 input of a BN layer under non-IID settings. Notably, we show that even if clients communicate
 63 after every local step, the *dependency of the gradient on the local mini-batch* prevents FEDAVG
 64 from recovering the gradient computed in the centralized training setting. We note that this does not
 65 happen to DNNs with GN, as GN does not use mini-batch statistics to normalize features. Taking
 66 this insight into account, we propose a simple yet highly effective treatment named FIXBN, which
 67 requires no architecture change, no additional training, and no extra communication costs.

68 2 Related Work

69 **Normalization layers in centralized training.** The benefits of BN [12] have been extensively studied
 70 in centralized training such as less internal covariate shift [12], smoother optimization landscape [32],
 71 robustness to hyperparameters [31] and initialization [35], accelerating convergence [12], etc. The
 72 noise of the estimated statistics of BN in mini-batch training is considered a regularizer [33] that
 73 improves generalization [12]. A recent study [39] shows that BN is still the irreplaceable normalizer
 74 vs. a wide range of alternative choices in general settings. Note that, unlike in FL, BN *often*
 75 *outperforms* GN in *standard centralized training*.

76 **Existing use of normalizers in FL.** In the context of FL, [14] is the first to suggest replacing BN with
 77 GN for non-IID decentralized learning. Several works [15, 40] report that LN can be competitive
 78 to GN. [41] enhances adversarial robustness by using statistics from reliable clients but not for
 79 improving performance. HETEROFL [42] proposed to simply normalize batch activations instead of
 80 tracking running statistics for the scenario that the clients have heterogeneous model architectures.
 81 These works aim to *replace* BN while *we analyze* BN and *reclaim its superiority*.

82 Several works [43, 44] propose dedicated server aggregation methods for BN statistics (separated
 83 from other model parameters) for specific tasks. For multi-modal learning, [45] proposes to maintain
 84 each modality as a different BN layer instead of sharing a single one. In personalized FL, [46, 47, 48]
 85 propose to maintain each client’s independent BN layer, inspired by the practice of domain adaptation
 86 in centralized training [49]. [50] leverages BN statistics to guide aggregation for personalization. The
 87 goals of these works are orthogonal to ours.

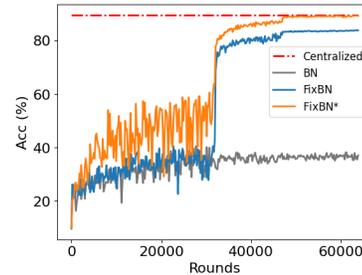


Figure 1: Our simple two-stage treatment FIXBN largely bridges the gaps of using BN in FL and centralized learning. Please see section 3 and subsection C.4 for more details about FIXBN (*: with SGD momentum) and this non-IID CIFAR-10 experiments.

88 3 Rethinking Batch Normalization in FL

89 3.1 Background

90 **Batch Normalization (BN).** The BN layer is widely used as a building block in feed-forward DNNs.
 91 Given an input feature vector \mathbf{h} , the BN layer normalizes the feature (via the mean $\mu_{\mathcal{B}}$ and variance
 92 $\sigma_{\mathcal{B}}^2$ computed on a batch of features \mathcal{B}), followed by a learnable affine transformation (via γ, β):
 93 $\hat{\mathbf{h}} = f_{\text{BN}}(\mathbf{h}; (\gamma, \beta), (\mu_{\mathcal{B}}, \sigma_{\mathcal{B}}^2)) = \gamma \frac{\mathbf{h} - \mu_{\mathcal{B}}}{\sqrt{\sigma_{\mathcal{B}}^2 + \epsilon}} + \beta$. ϵ is a small constant. In standard training, the
 94 statistics $\mu_{\mathcal{B}}$ and $\sigma_{\mathcal{B}}^2$ are computed on each training mini-batch during the forward passes. These
 95 mini-batch statistics are accumulated during training by the following exponential moving average
 96 (controlled by α) to replace $\mu_{\mathcal{B}}$ and $\sigma_{\mathcal{B}}^2$ in the equation above for testing:

$$\mu := \alpha\mu + (1 - \alpha)\mu_{\mathcal{B}}, \quad \sigma^2 := \alpha\sigma^2 + (1 - \alpha)\sigma_{\mathcal{B}}^2. \quad (1)$$

97 **Federated averaging (FEDAVG).** Federated learning (FL) learns a model parameterized by θ on
 98 the decentralized data $\mathcal{D}_m, \forall m \in [M]$ of M clients. For DNNs with BN layers θ includes learnable
 99 weights and the statistics $\{(\gamma, \beta), \mathbf{S}\}$ of all BN layers, where $\mathbf{S} = (\mu, \sigma^2)$ are the BN means and
 100 variances. The fundamental FL algorithm FEDAVG [2] solves it by multiple rounds of parallel local
 101 updates at the clients and global model aggregation at the server. Given an initial model $\bar{\theta}^{(0)}$, for
 102 round $t = 1, \dots, T$, FEDAVG performs:

$$\text{Local: } \theta_m^{(t)} = \text{ClientUpdate}(\mathcal{L}_m, \bar{\theta}^{(t-1)}); \quad \text{Global: } \bar{\theta}^{(t)} \leftarrow \sum_{m=1}^M \frac{|\mathcal{D}_m|}{|\mathcal{D}|} \theta_m^{(t)}. \quad (2)$$

103 During local training, the clients update the model parameters received from the server, typically by
 104 minimizing each client’s empirical risk \mathcal{L}_m with several steps (denoted as E) of mini-batch SGD. For
 105 the locally accumulated means and variances in BN, they are updated by Equation 1. During global
 106 aggregation, all the parameters in the locally updated models $\{\theta_m^{(t)}\}$, including the BN statistics, are
 107 averaged element-wise over clients. Typically, $E \gg 1$ due to communication constraints.

108 3.2 Problem: BN in FL cannot recover centralized performance

109 The non-IID issue is particularly problematic for DNNs with BN layers in FEDAVG since they depend
 110 on the activation mean and variance estimation computed on non-IID mini-batches.

111 We first consider communicating after every SGD step. That
 112 is, in the local training in Equation 2, we only perform a
 113 single mini-batch SGD update in each round, i.e., $E = 1$.
 114 At first glance, this should recover mini-batch SGD in cen-
 115 tralized learning (e.g., training on multi-GPUs with local
 116 shuffling). However, as shown in Table 1 and Figure 1, even
 117 with high-frequency communication after every SGD step,
 118 there is a huge accuracy gap (about 45%) between centralized
 119 and federated learning for DNNs with BN. As a reference,
 120 *such a gap very much disappears for DNNs with GN*. Intrigued by this observation, we investigate
 121 the potential reasons from three aspects below, focusing on the non-IID FL setting.

Table 1: **FL with communication every step** ($E = 1$). We train a ResNet20 with either BN or GN on the non-IID CIFAR-10 dataset (5 clients, 2 classes per client). Both the FL and centralized training use SGD without momentum.

	Norm	Centralized Acc.	FL Acc.
GN	87.46±0.57	87.37±1.16	
BN	89.30±0.89	42.93±2.75	

122 3.3 BN training dynamics

123 We first consider the properties of BN in standard training. We note that BN normalizes the activations
 124 in the forward pass to ensure stable forward and backward propagation [39]. A naive workaround
 125 for the non-IID issue is to force all clients to normalize with the same statistics. We investigate this
 126 idea by “decoupling” the updates of the model weights and the BN statistics. Specifically, under the
 127 high-frequency communication setting with $E = 1$, we modify Equation 2 as follows. (a) At round t ,
 128 given frozen weights in $\bar{\theta}^{(t)}$, we update local statistics $\{\mathbf{S}_m^{(t+1)}\}_{m=1}^M$ via Equation 1 and aggregate
 129 them into $\bar{\mathbf{S}}^{(t+1)}$. (b) We then locally update the model weights in the evaluation mode, using the
 130 global statistics $\bar{\mathbf{S}}^{(t+1)}$ to normalize the activations. (c) Finally, we aggregate the local models into
 131 $\bar{\theta}^{(t+1)}$. In the same FL experiment of subsection 3.2, we observe it achieves 52% accuracy, still far
 132 from the BN centralized performance 89%.

133 We hypothesize the weights and statistics need to collaborate carefully to enjoy the benefits of
 134 BN dynamics. First, using fixed statistics in local training sacrifices the “sampling” noise of the
 135 estimated statistics from different mini-batch $\mathcal{S}_B = (\boldsymbol{\mu}_B, \sigma_B^2)$, which is believed to help search a
 136 flatter loss landscape [33]. Second, using fixed statistics cannot properly normalize the activations in
 137 a mini-batch and could make DNN training fragile due to gradient explosion, especially in the earlier
 138 rounds of FEDAVG when the model weights and intermediate activations are changing rapidly.

139 3.4 Re-examining the BN statistics mismatch between training and testing

140 The reason why BN degrades in FL is believed to be the *statistics mismatch* issue pointed out
 141 by [14]. In section 5 of [14], the authors argued that since the local accumulated statistics
 142 $\{\mathcal{S}_m = (\boldsymbol{\mu}_{\mathcal{D}_m}, \sigma_{\mathcal{D}_m}^2)\}$ are estimated on each of the non-IID local data $\{\mathcal{D}_m\}$, their average could
 143 be significantly different from the true statistics of the global data $\mathcal{D} = \cup_m \mathcal{D}_m$. In other words, the
 144 average of $\{\mathcal{S}_m\}$ (over m) may not be ideal in testing. To verify its impact on performance, we
 145 design a simple experiment (details in Appendix B) aiming to *remove* the statistics mismatch.

146 After the entire FEDAVG is finished, we re-accumulate the statistics $\{(\boldsymbol{\mu}, \sigma^2)\}$ per BN layer directly
 147 on the aggregated training data $\mathcal{D} = \cup_m \mathcal{D}_m$ over clients, using Equation 1. Interestingly, we see a
 148 fairly small gain, i.e., less than 1% accuracy gain even on extreme non-IID settings.

149 *Based on the verification, we surmise that while the statistics mismatch problem indeed has a minor*
 150 *impact, it seems unlikely to account for the primary performance drops of BN in FL.*

151 3.5 BN makes the gradients biased in local training

152 We hypothesize that under the non-IID settings, the major reason for the performance drop comes from
 153 BN’s influence on local model training. As a simple illustration, we derive the forward-backward pass
 154 of the plain BN layer f_{BN} (see BN equation in subsection 3.1) for one example \mathbf{x}_i in a mini-batch \mathcal{B} .

$$\text{Forward: } \ell(\hat{\mathbf{x}}_i) = \ell(f_{\text{BN}}(\mathbf{x}_i; (\boldsymbol{\gamma}, \boldsymbol{\beta}), (\boldsymbol{\mu}_B, \sigma_B^2))) = \ell(\boldsymbol{\gamma} \frac{\mathbf{x}_i - \boldsymbol{\mu}_B}{\sqrt{\sigma_B^2 + \epsilon}} + \boldsymbol{\beta}) = \ell(\boldsymbol{\gamma} \tilde{\mathbf{x}}_i + \boldsymbol{\beta}); \quad (3)$$

$$\text{Backward: } \frac{\partial \ell}{\partial \mathbf{x}_i} = \frac{|\mathcal{B}| \frac{\partial \ell}{\partial \tilde{\mathbf{x}}_i} - \sum_{j=1}^{|\mathcal{B}|} \frac{\partial \ell}{\partial \tilde{\mathbf{x}}_j} - \tilde{\mathbf{x}}_i \cdot \sum_{j=1}^{|\mathcal{B}|} \frac{\partial \ell}{\partial \tilde{\mathbf{x}}_j} \cdot \tilde{\mathbf{x}}_j}{|\mathcal{B}| \sqrt{\sigma_B^2 + \epsilon}}, \quad (4)$$

155 where ℓ is an arbitrary loss function on the BN layer’s output $\hat{\mathbf{x}}_i$, “ \cdot ” is element-wise multiplication,
 156 and $\frac{\partial \ell}{\partial \tilde{\mathbf{x}}} = \boldsymbol{\gamma} \frac{\partial \ell}{\partial \hat{\mathbf{x}}}$. Please see Section 3 of [12] for the derivation of the gradient.

157 We can see that many terms in Equation 4 (colored in red) depend on the mini-batch features $\{\mathbf{x}_j\}_{j=1}^{|\mathcal{B}|}$
 158 or statistics $(\boldsymbol{\mu}_B, \sigma_B^2)$. The background gradient $\frac{\partial \ell}{\partial \mathbf{x}_i}$ w.r.t. the input vector \mathbf{x}_i is thus sensitive to
 159 what other examples in the mini-batch are. This is particularly problematic in FL on DNNs when
 160 clients’ data are non-IID. Suppose \mathbf{x}_i belongs to client m , the gradient $\frac{\partial \ell}{\partial \mathbf{x}_i}$ will be different when it
 161 is calculated locally with other data sampled only from \mathcal{D}_m and when it is calculated globally (in
 162 centralized training) with other data sampled from $\mathcal{D} = \cup_m \mathcal{D}_m$. Such bias will propagate to the latter
 163 layers in a DNN. Namely, even if communicating after every mini-match SGD step, *how a particular*
 164 *data example influences the DNN parameters is already different between FL and centralized training.*

165 *We surmise that this is the fundamental reason why DNNs with BN degrade in FL. Although it*
 166 *becomes quite intuitive after our elaboration, to our surprise, such a gradient issue was not explicitly*
 167 *pointed out by previous works¹. They mostly referred to the mismatch problem in [14].*

168 4 FIXBN: Towards a Proper Use of BN in Federated Learning

169 4.1 On fixing BN in FL

170 Given the analysis in section 3, we ask, *Is there a way to bypass the issues of BN in FL to reclaim its*
 171 *superior performance in centralized training?* We start our exploration by taking a deeper look into

¹We recently noticed that a concurrent work [51] pointed out this finding as well. Nevertheless, our analysis and solution are quite different from theirs.

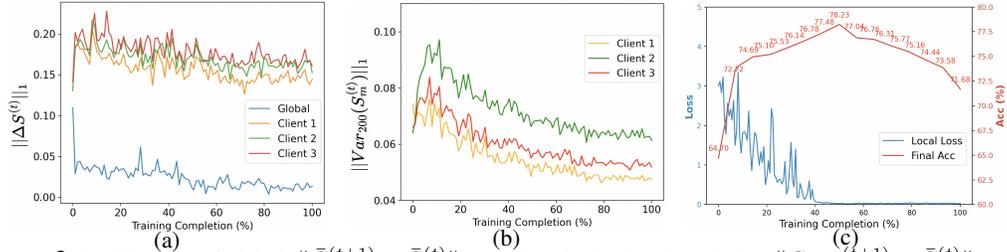


Figure 2: (a) Changes of global ($\|\bar{\mathcal{S}}^{(t+1)} - \bar{\mathcal{S}}^{(t)}\|_1$) and local mini-batch statistics ($\|\mathcal{S}_{m,B}^{(t+1)} - \bar{\mathcal{S}}^{(t)}\|_1$). (b) Variances (running over $t - 200$ to t) of local statistics $\mathcal{S}_m^{(t)}$. (c) Loss of global model on the training data and **final-round accuracy** when freezing BN statistics at different intermediate rounds (CIFAR-10, $E = 100$).

172 the dynamics of BN statistics during standard FEDAVG training. Under the same $E = 1$ experiments
 173 in subsection 3.2, we highlight two critical observations from Figure 2 (details in the captions).

174 First, as shown in Figure 2 (a), the local mini-batch statistics remain largely different from the global
 175 statistics, even at later rounds, which results from the discrepancy between non-IID local data. This is
 176 not surprising. However, it re-emphasizes the potentially huge impact of the issue in subsection 3.5.

177 Second, still in Figure 2 (a), we look at each curve alone. We find that both the global and local
 178 mini-batch statistics essentially converge. We further show the variances of the local statistics within
 179 each client become static in Figure 2 (b). This opens up the possibility to revisit the “decoupling”
 180 attempt in subsection 3.3.

181 Concretely, if local mini-batch statistics remain almost static in later rounds, replacing them with
 182 the fixed global statistics in local training may not degrade the benefits of BN. In contrast, it may
 183 fundamentally resolve the issue in subsection 3.5 — using the fixed global statistics in Equation 4
 184 could prevent local gradients from diverging. We investigate this idea by replacing local mini-batch
 185 statistics with fixed global statistics starting at different rounds. As shown in Figure 2 (c), if the round
 186 is chosen properly, the *final accuracy* can be largely improved. Based on this insight, we propose our
 187 FIXBN method to address the issues in section 3.

188 4.2 Two-stage training

189 To address the drawbacks discussed in section 3 simultaneously, we
 190 propose to divide FEDAVG training with BN into two stages, separated
 191 at round T^* , inspired by the widely-used decay learning strategy for
 192 SGD [52]. Supported by the insight in subsection 4.1, we first follow
 193 standard FEDAVG to explore a decent model solution space, thanks to
 194 BN’s important training dynamics as studied in subsection 3.3. Next, we
 195 propose to fine-tune the model for the rest of the training with the BN
 196 statistics *fixed*. This eliminates the statistics mismatch problem in sub-
 197 section 3.4 since now training and test share the same BN statistics. It
 198 also addresses the biased gradients caused by non-IID statistics in sub-
 199 section 3.5 as the local gradients no longer rely on mini-batches.

200 **Stage I: Exploration.** This stage is the standard FEDAVG with BN for
 201 two purposes: (a) to explore a proper model subspace without sacrificing
 202 BN’s benefits on optimization [12]; (b) to burn in the model and make it fitted to the training data. At
 203 the end, we save the aggregated statistics $\bar{\mathcal{S}}^{(T^*)}$ of each BN layer from the average model $\bar{\theta}^{(T^*)}$.

204 **Stage II: Calibration.** We anticipate that the exploration stage already finds a proper region of the
 205 model solution, and calibrated fine-tuning is needed to further improve the performance. Starting
 206 from round $T^* + 1$, we use the saved statistics as approximated global statistics to normalize the
 207 activations in local training. Since the model has been burned in, training with the fixed statistics is
 208 unlikely to suffer from the mentioned instability issue. In Figure 2 (c), we evaluate the training loss
 209 of the global model $\bar{\theta}^{(T^*)}$. It typically reaches a small loss after the first learning rate decay. Thus, *in*
 210 *experiments, we will simply fix the BN statistics since 50% of the total rounds of the FL training.*

211 While fairly simple, FIXBN cleverly leverages the global statistics to resolve the concerns in section 3,
 212 *with no architecture and objective change, no additional training and communication costs.*

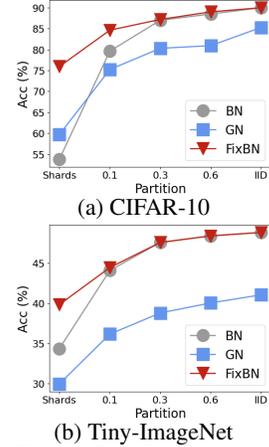


Figure 3: Non-IID partitions with $E = 100$ steps.

Table 4: **Comparison to other FL normalizer methods.** Test accuracy (%) of ResNet20 on CIFAR-10 given # of rounds. The setting follows [51].

FL Scheme	#R	IID	Non-IID
Centralized+BN	-	91.53	
Centralized+FIXBN	-	91.62	
FEDTAN [51]	580K	91.26	87.66
FEDAVG +BN	10K	91.35	45.96
FEDAVG +GN	10K	91.26	82.66
HETEROFL [42]	10K	91.21	30.62
FEDDNA [43]	10K	91.42	76.01
FEDAVG +FIXBN (Ours)	10K	91.35	87.71

Table 5: ResNet20 with different normalization layers FL on CIFAR-10 (Shards, $E = 100$).

Normalization Layer	Acc (%)
BN [12]	53.97 \pm 4.18
GN [16]	59.69 \pm 0.76
GN +WN [60]	66.90 \pm 0.81
LN [61]	54.54 \pm 1.21
IN [62]	59.76 \pm 0.43
FIXUP [63]	70.66 \pm 0.24
FIXBN (Ours)	76.56 \pm 0.66

213 5 Experiments (more in the appendix)

214 **Results on CIFAR-10 and Tiny-ImageNet.** We experiment on the FL benchmark CIFAR-10 [53]
 215 and Tiny-ImageNet [54] with 5/10 clients, ResNet20/ResNet18 [13], respectively. For the hyperpa-
 216 rameters, we generally follow [14] to use the SGD optimizer with 0.9 momentum, learning rate 0.02
 217 (decay by 0.1 at 50% and 75% of the total rounds, respectively), batch size 20, and full participation,
 218 $E = 100$. We train **fixed 128 epochs** of total local SGD updates over all the clients and communica-
 219 tion rounds. We consider **different non-IID degrees** including IID, Dirichlet(0.1, 0.3, 0.6) sampling
 220 follows [55], and Shards that each client only has data for 20% of the classes. We show FIXBN
 221 consistently outperforms BN and GN in Figure 3, especially in severe non-IID cases. Surprisingly,
 222 we found BN can sometimes outperform GN. A fine-grained comparison is in Appendix B.

223 **Results on ImageNet.** We extend FIXBN to
 224 ImageNet-1K [57] dataset which is split into 100
 225 clients Dirichlet (0.1) non-IID over classes. We learn
 226 a ResNet18 with 10% randomly sampled clients per
 227 round, 20 batch size, 0.1 learning rate (decay by 0.1
 228 every 30% of the total rounds), 2 local epochs, and 64
 229 epochs in total. Results in Table 2 show that FIXBN
 230 also perform the best.

231 **Comparison on realistic non-IID Cityscape.** We
 232 further consider a natural non-IID setting on the im-
 233 age segmentation dataset Cityscape [58]. We make
 234 each ‘‘city’’ a client and train 100 FEDAVG rounds
 235 using DeepLab-v3+ [59]. More details are in the
 236 appendix. Results in Table 3 show that FIXBN’s ef-
 237 fectiveness is generalized to different architectures
 238 and vision tasks.

239 **Other FL baselines.** We reproduce Table 1 in FEDTAN [51] to compare to other BN variants FL
 240 methods in Table 4. We note that FEDTAN requires communication rounds linear to the numbers of
 241 BN layers L as $\Theta(3L + 1)$, which is much more expensive than FIXBN. HETEROFL [42] directly
 242 normalizes the activations, which cannot resolve the non-IID issue.

243 **Beyond BN & GN, is there any FL-friendly alternative?** So far, we mainly focus on BN. Here
 244 we further compare to other normalization layers in Table 5. FIXBN still outperforms others.
 245 Interestingly, the normalization-free FIXUP [63] initialization for residual networks² performs much
 246 better than GN, suggesting a new alternative in FL besides FIXBN.

247 6 More Discussions and Conclusion

248 We revisit the use of BN layers and its common alternative, GN, in non-IID federated deep learning
 249 and conduct an in-depth analysis. We dissect the issues of BN in FL and propose a simple yet
 250 highly effective treatment named FIXBN to bridge the performance gap between FL and centralized
 251 training. We hope our study provides the community with a good foundation for the full (theoretical)
 252 understanding of the effectiveness of BN towards training deeper models in FL.

²Another concurrent work [64] also reports improving by replacing BN with scaled weight normalization, similar to [60] in Table 5.

Table 2: Class-non-IID ImageNet.

Method	Network	Acc. Δ_{-BN}
GN	ResNet18 [13]	33.33 \pm 0.57
BN		48.30 \pm 1.21
FIXBN		52.43 \pm 0.68 (+4.1)

Table 3: Pixel-wise accuracy and mean IoU (%) of image segmentation on Cityscapes.

Method	Backbone	Mean IoU Δ_{-BN}
GN	MobileNet-v2 [56]	43.2 \pm 0.33
BN		48.9 \pm 0.36
FIXBN		54.0 \pm 0.29 (+5.1)
GN	ResNet18 [13]	47.8 \pm 0.30
BN		52.6 \pm 0.38
FIXBN		57.2 \pm 0.32 (+4.6)

253 References

- 254 [1] Bommasani, R., D. A. Hudson, E. Adeli, et al. On the opportunities and risks of foundation
255 models. *arXiv preprint arXiv:2108.07258*, 2021.
- 256 [2] McMahan, H. B., E. Moore, D. Ramage, et al. Communication-efficient learning of deep
257 networks from decentralized data. In *AISTATS*. 2017.
- 258 [3] Zhou, F., G. Cong. On the convergence properties of a k -step averaging stochastic gradient
259 descent algorithm for nonconvex optimization. *arXiv preprint arXiv:1708.01012*, 2017.
- 260 [4] Stich, S. U. Local sgd converges fast and communicates little. In *ICLR*. 2019.
- 261 [5] Haddadpour, F., M. Mahdavi. On the convergence of local descent methods in federated learning.
262 *arXiv preprint arXiv:1910.14425*, 2019.
- 263 [6] Li, X., K. Huang, W. Yang, et al. On the convergence of fedavg on non-iid data. In *ICLR*. 2020.
- 264 [7] Zhao, Y., M. Li, L. Lai, et al. Federated learning with non-iid data. *arXiv preprint*
265 *arXiv:1806.00582*, 2018.
- 266 [8] Li, T., A. K. Sahu, M. Zaheer, et al. Federated optimization in heterogeneous networks.
267 *Proceedings of Machine Learning and Systems*, 2:429–450, 2020.
- 268 [9] Karimireddy, S. P., S. Kale, M. Mohri, et al. Scaffold: Stochastic controlled averaging for
269 federated learning. In *ICML*. 2020.
- 270 [10] Acar, D. A. E., Y. Zhao, R. Matas, et al. Federated learning based on dynamic regularization. In
271 *ICLR*. 2021.
- 272 [11] Chen, H.-Y., W.-L. Chao. Fedbe: Making bayesian model ensemble applicable to federated
273 learning. In *ICLR*. 2021.
- 274 [12] Ioffe, S., C. Szegedy. Batch normalization: Accelerating deep network training by reducing
275 internal covariate shift. In *International conference on machine learning*, pages 448–456.
276 PMLR, 2015.
- 277 [13] He, K., X. Zhang, S. Ren, et al. Deep residual learning for image recognition. In *CVPR*. 2016.
- 278 [14] Hsieh, K., A. Phanishayee, O. Mutlu, et al. The non-iid data quagmire of decentralized machine
279 learning. In *International Conference on Machine Learning*, pages 4387–4398. PMLR, 2020.
- 280 [15] Du, Z., J. Sun, A. Li, et al. Rethinking normalization methods in federated learning. In
281 *Proceedings of the 3rd International Workshop on Distributed Machine Learning*, pages 16–22.
282 2022.
- 283 [16] Wu, Y., K. He. Group normalization. In *Proceedings of the European conference on computer*
284 *vision (ECCV)*, pages 3–19. 2018.
- 285 [17] Jin, J., J. Ren, Y. Zhou, et al. Accelerated federated learning with decoupled adaptive opti-
286 mization. In *International Conference on Machine Learning*, pages 10298–10322. PMLR,
287 2022.
- 288 [18] Charles, Z., Z. Garrett, Z. Huo, et al. On large-cohort training for federated learning. *Advances*
289 *in neural information processing systems*, 34:20461–20475, 2021.
- 290 [19] Lin, T., L. Kong, S. U. Stich, et al. Ensemble distillation for robust model fusion in federated
291 learning. In *NeurIPS*. 2020.
- 292 [20] Yuan, H., W. Morningstar, L. Ning, et al. What do we mean by generalization in federated
293 learning? In *ICLR*. 2021.
- 294 [21] Reddi, S., Z. Charles, M. Zaheer, et al. Adaptive federated optimization. In *ICLR*. 2020.
- 295 [22] Hyeon-Woo, N., M. Ye-Bin, T.-H. Oh. Fedpara: Low-rank hadamard product for
296 communication-efficient federated learning. *arXiv preprint arXiv:2108.06098*, 2021.
- 297 [23] Yu, F., W. Zhang, Z. Qin, et al. Fed2: Feature-aligned federated learning. In *Proceedings of*
298 *the 27th ACM SIGKDD conference on knowledge discovery & data mining*, pages 2066–2074.
299 2021.
- 300 [24] Hosseini, H., H. Park, S. Yun, et al. Federated learning of user verification models without
301 sharing embeddings. In *International Conference on Machine Learning*, pages 4328–4336.
302 PMLR, 2021.

- 303 [25] Mohamad, M., J. Neubert, J. S. Argayo. Fedos: using open-set learning to stabilize training in
304 federated learning. *arXiv preprint arXiv:2208.11512*, 2022.
- 305 [26] Tenison, I., S. A. Sreeramadas, V. Mugunthan, et al. Gradient masked averaging for federated
306 learning. *arXiv preprint arXiv:2201.11986*, 2022.
- 307 [27] Yang, S., H. Hwang, D. Kim, et al. Towards the practical utility of federated learning in the
308 medical domain. *arXiv preprint arXiv:2207.03075*, 2022.
- 309 [28] PyTorch. Pytorch hub. <https://pytorch.org/hub/>, 2023.
- 310 [29] ONNX. Open neural network exchange (onnx) model zoo. [https://github.com/onnx/
311 models](https://github.com/onnx/models), 2023.
- 312 [30] Li, Y., G. Yuan, Y. Wen, et al. Efficientformer: Vision transformers at mobilenet speed. *Advances
313 in Neural Information Processing Systems*, 35:12934–12949, 2022.
- 314 [31] Bjorck, N., C. P. Gomes, B. Selman, et al. Understanding batch normalization. *Advances in
315 neural information processing systems*, 31, 2018.
- 316 [32] Santurkar, S., D. Tsipras, A. Ilyas, et al. How does batch normalization help optimization?
317 *Advances in neural information processing systems*, 31, 2018.
- 318 [33] Luo, P., X. Wang, W. Shao, et al. Towards understanding regularization in batch normalization.
319 In *ICLR*. 2019.
- 320 [34] Garbin, C., X. Zhu, O. Marques. Dropout vs. batch normalization: an empirical study of their
321 impact to deep learning. *Multimedia Tools and Applications*, 79(19):12777–12815, 2020.
- 322 [35] Yang, G., J. Pennington, V. Rao, et al. A mean field theory of batch normalization. *arXiv
323 preprint arXiv:1902.08129*, 2019.
- 324 [36] Wang, Y., Y. Xu, Q. Shi, et al. Quantized federated learning under transmission delay and
325 outage constraints. *IEEE Journal on Selected Areas in Communications*, 40(1):323–341, 2021.
- 326 [37] Zheng, S., C. Shen, X. Chen. Design and analysis of uplink and downlink communications
327 for federated learning. *IEEE Journal on Selected Areas in Communications*, 39(7):2150–2167,
328 2020.
- 329 [38] Chai, Z., Y. Chen, A. Anwar, et al. Fedat: a high-performance and communication-efficient fed-
330 erated learning system with asynchronous tiers. In *Proceedings of the International Conference
331 for High Performance Computing, Networking, Storage and Analysis*, pages 1–16. 2021.
- 332 [39] Lubana, E. S., R. Dick, H. Tanaka. Beyond batchnorm: towards a unified understanding of
333 normalization in deep learning. *Advances in Neural Information Processing Systems*, 34:4778–
334 4791, 2021.
- 335 [40] Zhang, G., M. Beitollahi, A. Bie, et al. Normalization is all you need: Understanding layer-
336 normalized federated learning under extreme label shift. *arXiv preprint arXiv:2308.09565*,
337 2023.
- 338 [41] Hong, J., H. Wang, Z. Wang, et al. Federated robustness propagation: Sharing adversarial
339 robustness in federated learning. *arXiv preprint arXiv:2106.10196*, 1, 2021.
- 340 [42] Diao, E., J. Ding, V. Tarokh. Heterofl: Computation and communication efficient federated
341 learning for heterogeneous clients. In *ICLR*. 2020.
- 342 [43] Duan, J.-H., W. Li, S. Lu. Feddna: Federated learning with decoupled normalization-layer ag-
343 gregation for non-iid data. In *Joint European Conference on Machine Learning and Knowledge
344 Discovery in Databases*, pages 722–737. Springer, 2021.
- 345 [44] Idrissi, M. J., I. Berrada, G. Noubir. Fedbbs: Learning on non-iid data in federated learning
346 using batch normalization. In *2021 IEEE 33rd International Conference on Tools with Artificial
347 Intelligence (ICTAI)*, pages 861–867. IEEE, 2021.
- 348 [45] Bernecker, T., A. Peters, C. L. Schlett, et al. Fednorm: Modality-based normalization in
349 federated learning for multi-modal liver segmentation. *arXiv preprint arXiv:2205.11096*, 2022.
- 350 [46] Li, X., M. JIANG, X. Zhang, et al. Fed{bn}: Federated learning on non-{iid} features via local
351 batch normalization. In *ICLR*. 2021.
- 352 [47] Andreux, M., J. O. d. Terrail, C. Beguier, et al. Siloed federated learning for multi-centric
353 histopathology datasets. In *Domain Adaptation and Representation Transfer, and Distributed
354 and Collaborative Learning*, pages 129–139. Springer, 2020.

- 355 [48] Jiang, M., X. Zhang, M. Kamp, et al. Tsmobn: Interventional generalization for unseen clients
356 in federated learning. *arXiv preprint arXiv:2110.09974*, 2021.
- 357 [49] Li, Y., N. Wang, J. Shi, et al. Revisiting batch normalization for practical domain adaptation.
358 *arXiv preprint arXiv:1603.04779*, 2016.
- 359 [50] Lu, W., J. Wang, Y. Chen, et al. Personalized federated learning with adaptive batchnorm for
360 healthcare. *IEEE Transactions on Big Data*, 2022.
- 361 [51] Wang, Y., Q. Shi, T.-H. Chang. Why batch normalization damage federated learning on non-iid
362 data? *arXiv preprint arXiv:2301.02982*, 2023.
- 363 [52] Robbins, H., S. Monro. A stochastic approximation method. *The annals of mathematical
364 statistics*, pages 400–407, 1951.
- 365 [53] Krizhevsky, A., G. Hinton, et al. Learning multiple layers of features from tiny images. 2009.
- 366 [54] Le, Y., X. Yang. Tiny imagenet visual recognition challenge. *CS 231N*, 2015.
- 367 [55] Hsu, T.-M. H., H. Qi, M. Brown. Measuring the effects of non-identical data distribution for
368 federated visual classification. *arXiv preprint arXiv:1909.06335*, 2019.
- 369 [56] Sandler, M., A. Howard, M. Zhu, et al. Mobilenetv2: Inverted residuals and linear bottlenecks.
370 In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages
371 4510–4520. 2018.
- 372 [57] Deng, J., W. Dong, R. Socher, et al. Imagenet: A large-scale hierarchical image database. In
373 *2009 IEEE Conference on Computer Vision and Pattern Recognition*, pages 248–255. 2009.
- 374 [58] Cordts, M., M. Omran, S. Ramos, et al. The cityscapes dataset. In *CVPR Workshop on The
375 Future of Datasets in Vision*. 2015.
- 376 [59] Chen, L.-C., Y. Zhu, G. Papandreou, et al. Encoder-decoder with atrous separable convolution
377 for semantic image segmentation. In *Proceedings of the European conference on computer
378 vision (ECCV)*, pages 801–818. 2018.
- 379 [60] Qiao, S., H. Wang, C. Liu, et al. Micro-batch training with batch-channel normalization and
380 weight standardization. *arXiv preprint arXiv:1903.10520*, 2019.
- 381 [61] Ba, J. L., J. R. Kiros, G. E. Hinton. Layer normalization. *arXiv preprint arXiv:1607.06450*,
382 2016.
- 383 [62] Ulyanov, D., A. Vedaldi, V. Lempitsky. Instance normalization: The missing ingredient for fast
384 stylization. *arXiv preprint arXiv:1607.08022*, 2016.
- 385 [63] Zhang, H., Y. N. Dauphin, T. Ma. Fixup initialization: Residual learning without normalization.
386 *arXiv preprint arXiv:1901.09321*, 2019.
- 387 [64] Zhuang, W., L. Lyu. Is normalization indispensable for multi-domain federated learning? *arXiv
388 preprint arXiv:2306.05879*, 2023.
- 389 [65] LeCun, Y., L. Bottou, Y. Bengio, et al. Gradient-based learning applied to document recognition.
390 *Proceedings of the IEEE*, 86(11):2278–2324, 1998.
- 391 [66] Karimireddy, S. P., M. Jaggi, S. Kale, et al. Mime: Mimicking centralized stochastic algorithms
392 in federated learning. *arXiv preprint arXiv:2008.03606*, 2020.

Appendix

393

394 We provide details omitted in the main paper.

- 395 • Appendix A: details of experimental setups (cf. Appendix B and section 4 of the main
396 paper).
- 397 • Appendix B: experimental results and analysis for BN vs GN (cf. section 5 and section 4 of
398 the main paper).
- 399 • Appendix C: additional experimental results and analysis for BN parameters and other
400 ablation studies (cf. Appendix B and section 4 of the main paper).

Table F: Summary of datasets and setups.

Dataset	Task	#Class	#Training	#Test/Valid	#Clients	Resolution	Networks
CIFAR-10	Classification	10	50K	10K	5 ~ 100	32 ²	LeNet-CNN, ResNet-20
Tiny-ImageNet	Classification	200	100K	10K	10	64 ²	ResNet-18
ImageNet	Classification	1,000	1,200K	100K	100	224 ²	ResNet-18
Cityscapes	Segmentation	19	3K	0.5K	18	768 ²	DeepLabv3 + {MobileNet-v2, ResNet-50}

Table G: Default FL settings and training hyperparameters in the main paper.

Dataset	Non-IID	Sampling	Optimizer	Learning rate	Batch size	T^* for FIXBN
CIFAR-10	Shards, Dirichlet({0.1, 0.3, 0.6}), IID	10 ~ 100%	SGD + 0.9 momentum	0.2/0.02	20	50% of total rounds
Tiny-ImageNet	Shards, Dirichlet({0.1, 0.3, 0.6}), IID	50%	SGD + 0.9 momentum	0.02	20	50% of total rounds
ImageNet	Dirichlet 0.1	10%	SGD + 0.9 momentum	0.1	20	50% of total rounds
Cityscapes	Cities	50%	Adam	0.01/0.001	8	90th round

401 A Experiment Details

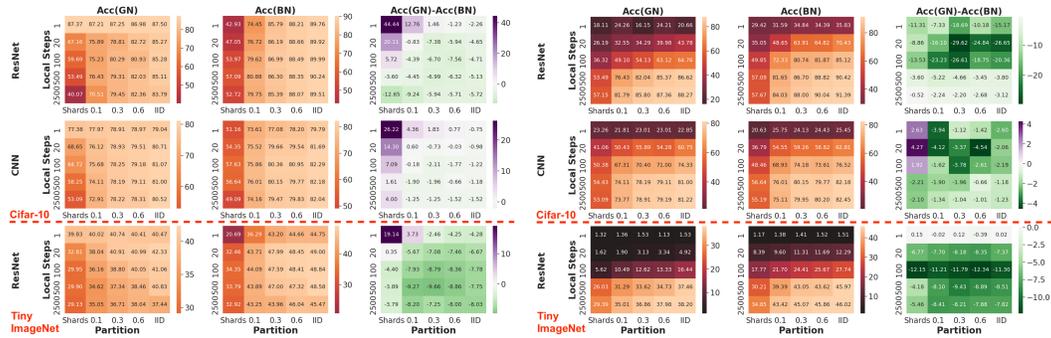
402 A.1 Datasets, FL settings, and hyperparameters

403 We use FEDAVG for our studies, with weight decay $1e-4$ for local training. Learning rates are
404 decayed by 0.1 at 50%, 75% of the total rounds, respectively. Besides that, we summarize the
405 training hyperparameters for each of the federated experiments included in the main paper in Table G.
406 Additionally, for the Cityscape experiments in Table 3, we make each “city” a client and run 100
407 rounds, with local steps to be 5 epochs. More details about the datasets are provided in Table F.

408 For pre-processing, we generally follow the standard practice which normalizes the images and
409 applies some augmentations. CIFAR-10 images are padded 2 pixels on each side, randomly flipped
410 horizontally, and then randomly cropped back to 32×32 . For Tiny-ImageNet, we simply randomly
411 cropped to the desired sizes and flipped horizontally following the official PyTorch ImageNet training
412 script. For the Cityscapes dataset, we use output stride 16. In training, the images are randomly
413 cropped to 768×768 and resized to 2048×1024 in testing.

414 B A Detailed Study of BN vs. GN

415 Results in section 5 come to an unexpected finding: BN outperforms GN in many cases, contradicting
416 the common belief that one should replace BN with GN in FL proposed in [14] and followed by



(a) Fixed 128 epochs (b) Fixed 128 rounds

Figure D: Is GN always better than BN? No. We compare their test accuracy in various FL settings on CIFAR-10 and Tiny-ImageNet, including different non-IID partitions and numbers of local steps E . Fixed budget of the total number of SGD updates (e.g., for CIFAR-10, $20 E \times 5$ clients \times 3200 rounds = 128 epochs) or the number of total rounds (128 rounds) are given.

many works summarized in section 1 and section 2. To answer this question, we revisit the study in [14] (which considers mere one FL setting) and provide a detailed study to compare BN and GN by varying several critical factors in FL to have a more complete picture.

Experiment setup. We focus on CIFAR-10 [53] and Tiny-ImageNet [54] datasets, following the setup in section 5. We consider more factors like (1) **degrees of non-IID**, ordered in increasing skewness: **IID, Dirichlet(0.1, 0.3, 0.6), and Shards**. As practical FL is constrained on computation, we consider two (2) **budget criteria**: **fixed 128 epochs** of total local SGD updates over all the clients and communication rounds, and **fixed 128 rounds** of communication. In every round, each client runs $\{1, 20, 100, 500, 2500\}$ of (3) **local steps (E)**. We further include LeNet-like CNN [65] for CIFAR-10.

B.1 Revisiting: Is GN really better than BN?

Observations. We highlight the following observation from Figure D, augmenting the findings in [14]:

- **No definite winners.** GN is often considered the default replacement for BN in previous FL works (section 1 and section 2). However, according to Figure D, GN is not always better than BN.
- **BN often outperforms GN.** Instead, in most settings, BN outperforms GN. This can be seen from the green cells in “Acc(GN)-Acc(BN)” heatmaps of Figure D.
- **GN outperforms BN merely in extreme cases.** We find that GN outperforms BN (the purple cells in “Acc(GN)-Acc(BN)” heatmaps) only in the extreme non-IID (e.g., Shards) and highly frequent communication (e.g., $E = 1$) settings. When clients cannot communicate frequently, the case where many existing FL works focus on, BN seems to be the better choice for normalization.
- **The trends along the number of local steps E per communication round.** It is a perhaps well-known fact that increasing the number of local steps leads to greater drift as the local models become more biased [9]. However, using more local steps also allows for more updates to the local models, potentially leading to an improved average model. To balance these competing considerations, we will discuss two criteria. For (a) **fixed epochs over all communication rounds**, a larger number of local steps means fewer communication rounds, in which GN degrades monotonically “as expected”. *Interestingly*, BN has an opposite trend. BN actually improves and outperforms GN with larger E s. For (b) **fixed rounds**, understandably, using more local steps improves both GN and BN, since more local SGD updates are made in total. Nevertheless, the improvement saturates (e.g., $E \geq 500$).
- **Small difference from statistics mismatch.** In subsection 3.4, we discuss that the BN statistics mismatch problem might be minor. We re-estimate the statistics on global data and see a negligible accuracy gain from 44.09% to 44.87% on the Tiny-ImageNet (Dir(0.1), fixed epochs, $E = 100$).
- **More settings.** We verify in the next section that factors like participation rates and the number of clients for partitioning the data do not change the above observation.

Additional figures. At the beginning of this section, we provide a detailed empirical study to compare BN and GN across various FL settings to understand their sweet spots. Here we provide a closer look at the observations we summarized above.

- **The trends along the number of local steps E per communication round.** In subsection B.2, we identify the opposite trends along #local steps E between BN and GN. As shown in Figure F, we see GN drops with less communication as expected due to the well-known non-IID model drift problem in FL. Interestingly, we found that BN can actually improve within a certain range of communication frequencies (for local steps in $[1,500]$), which suggests that further investigation and theoretical analysis are required for BN in FL.
- **More settings.** We further verify that factors such as participation rate and the number of clients for partitioning the data in Figure G. As expected, the results are consistent with the observations summarized in subsection B.1, particularly in that there is no definite winner between BN and GN while BN often outperforms GN.

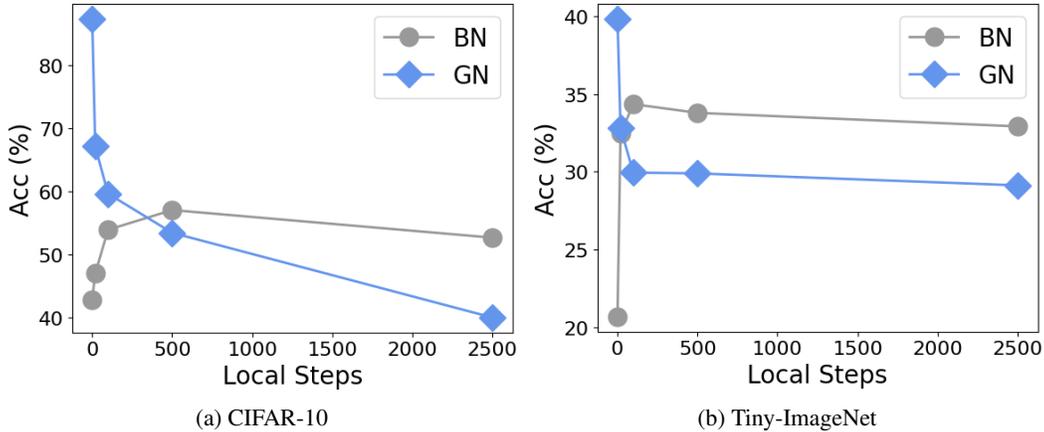


Figure F: **The opposite trends along #local steps E .** We consider the (Shards, **fixed epochs**) setting: the more the local step E is, the fewer the total number of communication rounds is. GN drops with less communication as expected, while BN can improve.

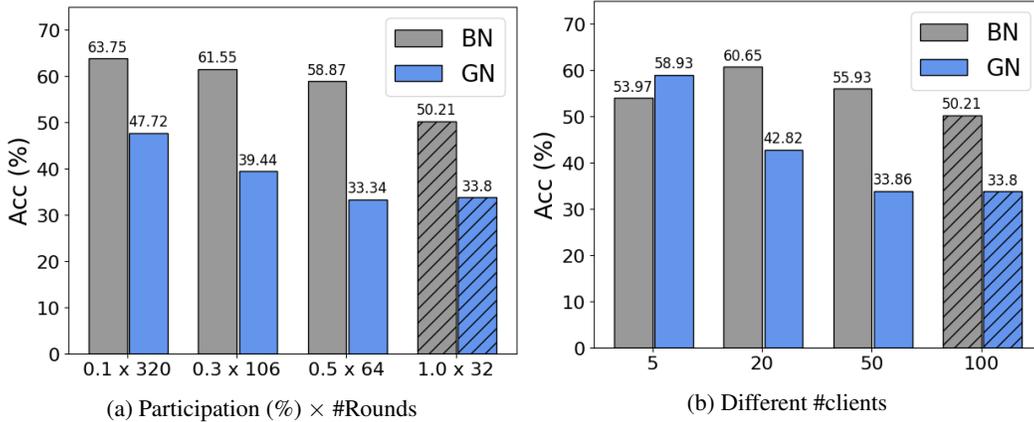


Figure G: **More settings.** We consider more clients ($M = 5 \sim 100$, $E = 100$) for partitioning CIFAR-10 (Shards) with fixed epochs and varying the participation rate of clients every round.

466 B.2 Effects of communication frequency

467 The constraint in communication, i.e., clients cannot aggregate the gradients frequently, is commonly
 468 believed as a major reason that hinders the performance of FL due to model drift [9]. As BN cannot
 469 recover the centralized gradient even with high communication frequency and is outperformed by GN
 470 in such a setting, one may expect that BN will be consistently surpassed by GN when the frequency
 471 drops. But surprisingly, as observed in subsection B.1, BN is *unreasonably effective when training*
 472 *with fewer communication rounds but more local steps per round.*

473 In Figure E, we vary the number of local SGD steps per communi-
 474 cation round (i.e., E) but fix the total number of SGD steps. We see
 475 the drastically different effect of E on BN and GN. In particular,
 476 while the performance of GN drops along with increasing E , BN
 477 somehow benefits from a larger E . Such a discrepancy suggests the
 478 need for a deeper (theoretical) analysis of the usage of BN in FL.

479 C Additional experimental results and analysis

480 C.1 Additional study of fixing BN parameters

481 In subsection 3.4, we discuss that the BN statistics are the main critical
 482 parameters in FL and thus motivate our design in FIXBN to fix
 483 the BN statistics to be the global aggregated ones after certain rounds.
 484 Here we include a further study to confirm the importance of BN
 485 statistics by comparing them with the learnable affine transformation
 486 parameterized by (γ, β) .

487 For FIXBN, besides fixing the BN statistics at round T^* , we consider
 488 fixing the (γ, β) alone or together. The results on CIFAR-10 (Shards,
 489 fixed epochs, $E = 100$) setting using ResNet20 is in Table H. We
 490 observe that fixing the (γ, β) only has slight effects on the test accuracy either in combination with
 491 fixing (γ, β) or not, validating that the statistics are the main reason making it suffers more in FL,
 492 compared to the affine transformation. Fixing (γ, β) alone cannot match the performance of the
 493 originally proposed FIXBN.

Table H: **Fixing different parameters as in FIXBN.** We consider fixing the BN statistics (μ, σ) as in original FIXBN or fixing the parameters (γ, β) of the affine transformation in BN layers. on CIFAR-10 (Shards, fixed epochs, $E = 100$) setting using ResNet20.

(μ, σ)	(γ, β)	Acc (%)
✓	✓	75.22
✓	✗	76.56
✗	✓	55.33
✗	✗	53.97

494 C.2 Different # of groups for GN

495 For experiments in our study, we set the # of groups = 2 for GN layers. We did not find the group
 496 size a significant factor for the performance, as confirmed in Table I.

Table I: **Effects of the groupsize for GN.** We experiment with different # of groups (2 ~ 8) to divide the channels in GN layers in the CIFAR-10 (Shards, $E = 100$) with fixed epochs setting.

Groupsize	Acc(%)
2	59.42
4	57.61
8	58.86

497 C.3 Effects of Batch Size for BN

498 We experiment with various batch sizes for both BN and FIXBN in the CIFAR-10 (Shards, $E = 1$)
 499 setting and saw FIXBN maintains the advantage over standard FEDAVG +BN.

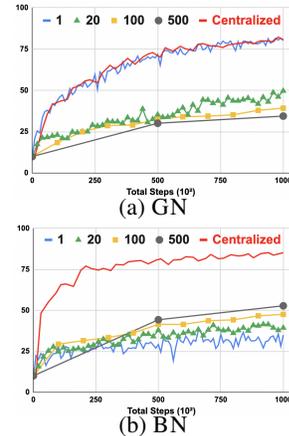


Figure E: Test accuracy on CIFAR-10 for different local steps (E) per communication given a fixed number of SGD steps.

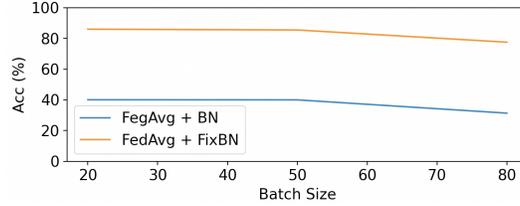


Figure H: FIXBN maintains advantage over different batch size selections.

500 C.4 Maintained SGD momentum further bridge FL to centralized performance

501 **Maintained SGD momentum.** Besides BN, we identify another gap between FEDAVG and centralized training. While using SGD momentum in standard FEDAVG during local training is common, it will be discarded at the end of the round and re-initialized (along with any optimizer states) at the beginning of the next round of local training in FEDAVG. That is, the first several SGD steps in a round cannot benefit from it.

509 To further bridge the gap, we present a fairly simple method, which is to keep the **local momentum** without re-initialization after the end of the local training in each round. This makes it a stateful method suitable for cross-silo FL. Another stateless choice is to maintain **global momentum** [66] by uploading the local momentum to the server in every round and aggregating it with Equation 2, for initializing the momentum of the next round of local training, with the cost of double message size. Empirically, we found the two methods yield similar gains (as will be shown in Figure I) and recover centralized performance if communicating every step (Figure 1).

522 **Experimental Results.** We combine each normalizer with the maintained local momentum and global momentum proposed in subsection C.4, respectively. We show FIXBN’s effectiveness against BN and GN in Figure I in the (Shards, fixed epoch) setting with different numbers of local steps per communication E of $\{1, 20, 100, 500, 2500\}$. We see FIXBN performs consistently better. More importantly, FIXBN remains highly accurate in fast communication, unlike BN, confirming that it mitigates the deviation issue in subsection 3.5 well. The improvements of using maintained global/local momentum are similar, providing the flexibility of stateless/stateful use cases. More gains are at small E , supporting our motivation to fix the zero initialization issue of the momentum. Across different settings, we see $\text{FixBN} \geq \text{BN} > \text{GN}$ in performance, consistent with Figure 3.

533 Both of them improve BN notably, especially at small E , supporting our motivation to fix the zero initialization issue of the local SGD momentum to stabilize the gradients. Indeed, in Figure 1 with $E = 1$, we show FIXBN largely recovers centralized performance, making BN much more applicable in FL.

537 C.5 Training curves

538 We provide the training curves of FIXBN and other normalizers under various settings in fixed 128 epochs using ResNet20 in Figure J, Figure K, Figure L, and Figure M, corresponding to Appendix B.

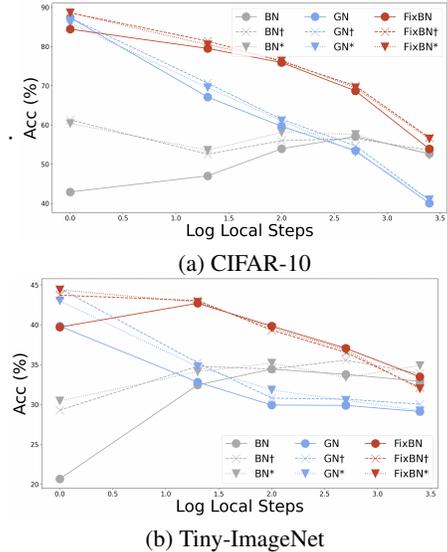


Figure I: **Maintained momentum.** Normalizers augmented with maintained **global momentum** (\dagger) and **local momentum** ($*$) with different numbers of local steps per communication E .

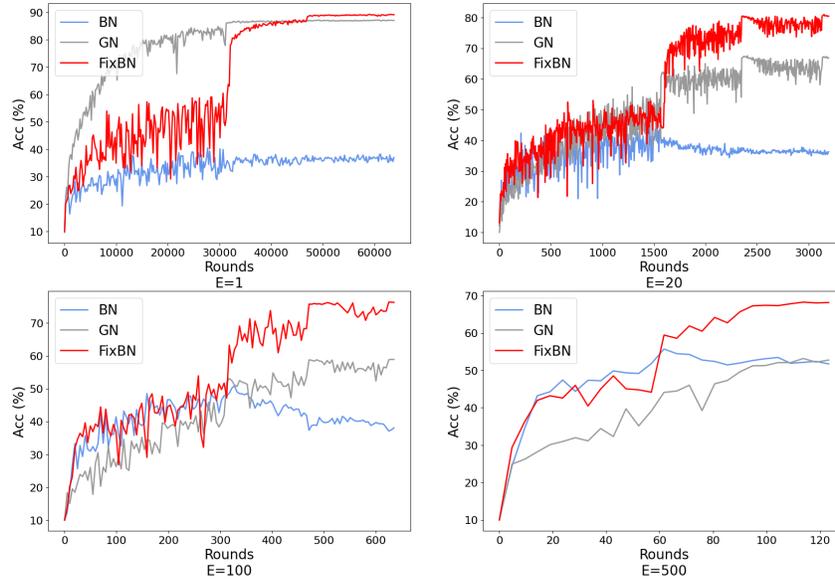


Figure J: Convergence curves of the test accuracy of CIFAR-10 with fixed epoch and **Shards non-IID** partitions, with $E = 1 \sim 500$.

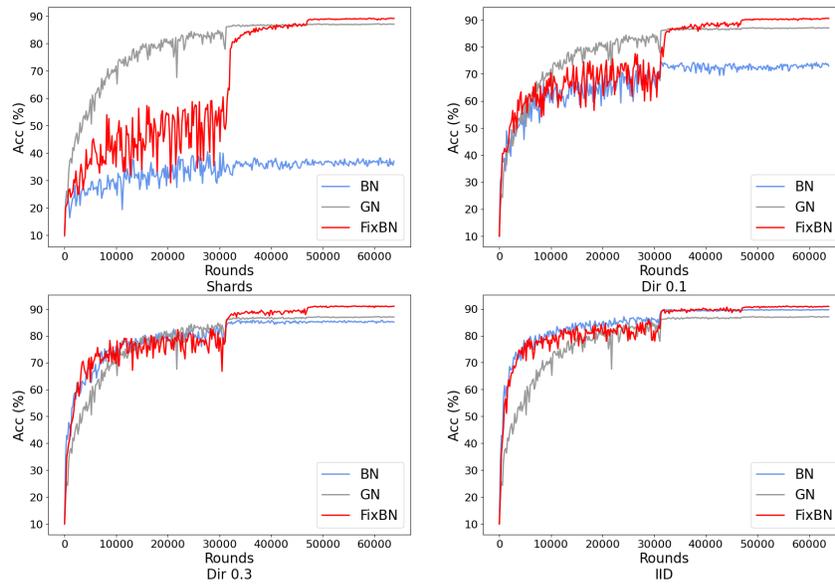


Figure K: Convergence curves of the test accuracy of CIFAR-10 in fixed epoch, **different non-IID partitions**, and $E = 1$ setting.

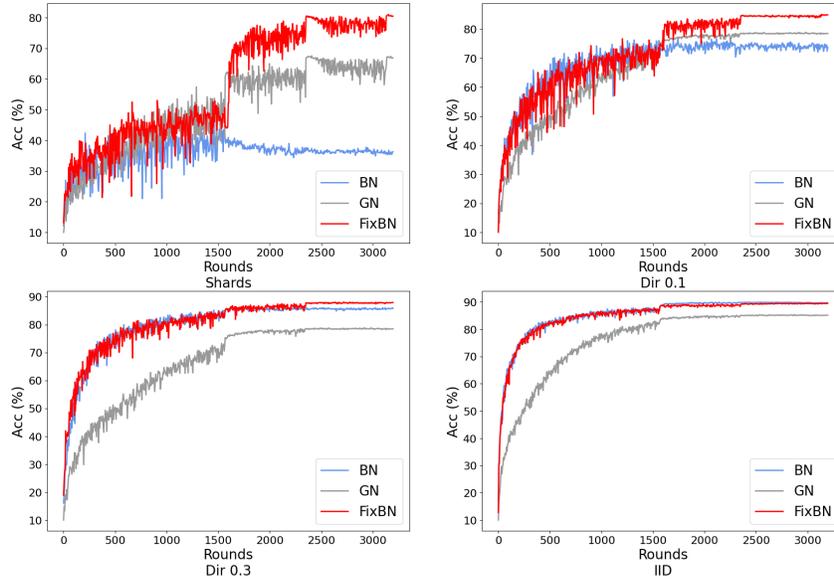


Figure L: Convergence curves of the test accuracy of CIFAR-10 in fixed epoch, **different non-IID partitions**, and $E = 20$ setting.

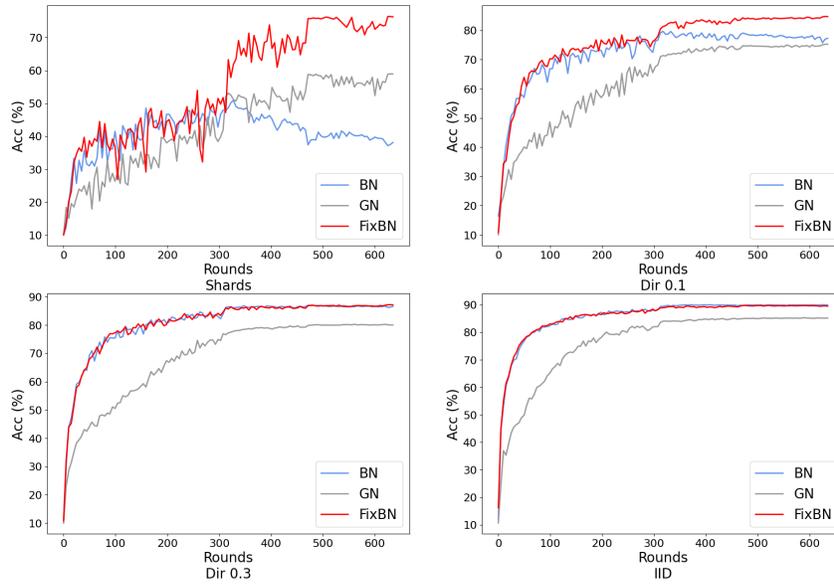


Figure M: Convergence curves of the test accuracy of CIFAR-10 in fixed epoch, **different non-IID partitions**, and $E = 100$ setting.