DISTRIBUTED QUASI-NEWTON METHOD FOR FAIR AND FAST FEDERATED LEARNING

Anonymous authors

Paper under double-blind review

ABSTRACT

Federated Learning (FL) is a promising technology that enables edge devices/clients to collaboratively and iteratively train a machine learning model under the coordination of a central server. The most common approach to FL is first-order methods, where clients send their local gradients to the server in each iteration. However, these methods often suffer from slow convergence rates. As a remedy, secondorder methods, such as quasi-Newton, can be employed in FL to accelerate its convergence. Unfortunately, similarly to the first-order FL methods, the application of second-order methods in FL can lead to unfair models, achieving high average accuracy while performing poorly on certain clients' local datasets. To tackle this issue, in this paper we introduce a novel second-order FL framework, dubbed distributed quasi-Newton federated learning (DQN-Fed). This approach seeks to ensure fairness while leveraging the fast convergence properties of quasi-Newton methods in the FL context. Specifically, DQN-Fed helps the server update the global model in such a way that (i) all local loss functions decrease to promote fairness, and (ii) the rate of change in local loss functions aligns with that of the quasi-Newton method. We prove the convergence of DQN-Fed and demonstrate its *linear-quadratic* convergence rate. Moreover, we validate the efficacy of DON-Fed across a range of federated datasets, showing that it surpasses state-of-the-art fair FL methods in fairness, average accuracy and convergence speed. The Code for paper is publicly available at https://github.com/ICMLDQNFed/ICMLDQN.

029 030 031

032

004

006

008 009

010 011

012

013

014

015

016

017

018

019

021

024

025

026

027

028

1 INTRODUCTION

Traditionally, machine learning (ML) models are trained centrally, with data stored in a central server. However, in modern applications, devices often resist sharing private data remotely. To address this, federated learning (FL) was introduced by McMahan et al. (2017), where each device trains locally with a central server. In FL, devices share only local updates, maintaining data privacy. FedAvg, proposed by McMahan et al. (2017), is a popular first-order FL method. It combines local stochastic gradient descent (SGD) on each client with iterative model averaging. The server sends the global model to selected clients Eichner et al. (2019); Wang et al. (2021a), which perform local SGD on their training data. Local gradients are sent back to the server, which calculates their (weighted) average to update the global model iteratively.

Nevertheless, first-order FL methods tend to exhibit slow convergence, particularly in terms of the number of iterations or communication rounds required (Krouka et al., 2022). More precisely, the convergence rate of first-order FL algorithms is sublinear, i.e., the required number of communication rounds T_{ϵ} to achieve ϵ -accurate solution is $T_{\epsilon} = \mathcal{O}(\frac{1}{\epsilon})$. Additionally, their convergence speed is highly influenced by the condition number, which is dependent on several factors, including: (i) the architecture of the model being trained, (ii) the choice of loss function, and (iii) the distribution of the training data (Elgabli et al., 2022).

To overcome this limitation, second-order methods can be applied in FL to significantly boost convergence speed (Safaryan et al., 2022; Elgabli et al., 2022). By estimating the local curvature of the loss landscape, these methods provide more adaptive and efficient update directions, leading to faster and more reliable convergence (Battiti, 1992). Specifically, in second-order FL methods, the clients compute the Newton direction for their respective local loss functions and send these directions to the server. The server then averages the Newton directions from all clients and updates

the global model in the direction of this average (Ghosh et al., 2020; Zhang & Lin, 2015). Moreover, since Newton's methods require calculating the inverse of the Hessian matrix at each iteration—a computationally expensive operation—the inverse is typically approximated using iterative techniques, leading to quasi-Newton methods (Wang et al., 2018).

While Newton-type methods accelerate the convergence of FL algorithms, they do not guarantee that the averaged Newton direction computed by the server is a descent direction for all clients—a limitation also present in first-order FL methods (Hu et al., 2022; Pan et al., 2024; Chen et al., 2024). In other words, upon updating the global model toward this averaged direction, the loss function for some client my not decrease, potentially leading to poor performance on their private datasets. As a result, the learned model might exhibit *unfairness*, with high average accuracy but poor performance for clients whose data distributions differ from the majority¹. Thus, naively applying Newton methods in FL can lead to the training of an *unfair* model (see Section 5 for results).

066 To tackle the issue mentioned above, this paper presents a novel second-order FL framework, dubbed 067 distributed quasi-Newton federated learning (DQN-Fed). This approach aims to ensure fairness while 068 leveraging the fast convergence properties of quasi-Newton methods in the FL setting. In particular, 069 DQN-Fed is designed to assist the server in updating the global model such that (i) all local loss functions decrease resulting in training a fair model, and (ii) the rate of change in local loss functions 071 aligns with the rate of change in the quasi-Newton method. To achieve this, based on the received local quasi-Newton directions and the local gradients, the server identifies an updating direction that 072 satisfies both of the aforementioned conditions. This will in turn yield a fair FL algorithm, as the 073 global updating direction is descent for all the clients. Moreover, the convergence of DQN-Fed is 074 fast, as the rate of change in the local loss functions follow quasi-Newton methods. 075

In summary, the contributions of the paper are as follows:

• We introduce distributed quasi-Newton federated learning (DQN-Fed), a method designed to assist the server in updating the global model to achieve both fairness and fast convergence in FL.

• We present a closed-form solution for calculating the global updating direction, distinguishing our approach from many existing fair FL methods that depend on iterative or generic quadratic programming techniques.

• Leveraging common assumptions in FL literature, we establish the convergence proof for DQN-Fed algorithm across various FL setups. In addition, we prove the convergence rate of the proposed method, and show that DQN-Fed exhibits a linear-quadratic convergence rate. Specifically, the convergence is either quadratic, with $T_{\epsilon} = \mathcal{O}\left(\log \log \frac{1}{\epsilon}\right)$, or linear, with $T_{\epsilon} = \mathcal{O}\left(\frac{1}{\log(\frac{\lambda}{L\delta})}\log \frac{1}{\epsilon}\right)$, where λ , L and δ are constants.

• Through *comprehensive* experiments conducted on <u>seven</u> different datasets (six vision datasets and one language dataset), we demonstrate that DQN-Fed attains superior fairness level among clients, and converges faster compared to the state-of-the-art fair alternatives.

2 RELATED WORKS

093 094 095

096

098

100

101

102

103

104

089

090

091 092

077

078

• Fairness in FL. The literature offers a myriad of perspectives to address the challenge of fairness in FL. These methods include client selection Nishio & Yonetani (2019); Huang et al. (2020a; 2022); Yang et al. (2021), contribution Evaluation Zhang et al. (2020); Lyu et al. (2020); Song et al. (2021); Le et al. (2021), incentive mechanisms Zhang et al. (2021); Kang et al. (2019); Ye et al. (2020); Zhang et al. (2020), and the methods based on the loss function. Specifically, our work falls into the latter category. This approach aims to achieve uniform test accuracy across clients. In particular, works within this framework focus on reducing the variance of test accuracy among participating clients. We provide a thorough review on fairness issue in ML and FL in Appendix L.

• Second-Order FL methods. DistributedNewton (Ghosh et al., 2020) and LocalNewton (Gupta et al., 2021) perform Newton's method instead of SGD on local machines to accelerate the convergence of local models. FedNew (Elgabli et al., 2022) utilized one pass ADMM on local machines

¹⁰⁵ 106 107

¹Learning an unfair model is a common challenge in first-order FL methods as well, and there is a substantial body of research dedicated to developing fair FL models (Mohri et al., 2019; Du et al., 2021; Li et al., 2020; Hu et al., 2022; Hamidi & YANG, 2024).

to calculating local directions and approximate Newton method to update the global model. FedNL (Safaryan et al., 2022) send the compressed local Hessian updates to global server and performed
Newton step globally. Based on eigendecomposition of the local Hessian matrices, SHED (Dal Fabbro et al., 2024) incrementally updated eigenvector-eigenvalue pairs to the global server and recovered
the Hessian to use Newton method. Recently, Li et al. (2024) proposed federated Newton sketch
methods (FedNS) to approximate the centralized Newton's method by communicating the sketched
square-root Hessian instead of the exact Hessian.

115 116

117 118

119

125

127

128

129

132

133 134 135

146

3 NOTATION AND PRELIMINARIES

3.1 NOTATION

120 We denote by [K] the set of integers $\{1, 2, \dots, K\}$. In addition, we define $\{f_k\}_{k \in [K]} = \{f_1, f_2, \dots, f_K\}$ for a scalar/function f. We use bold small letters to represent vectors, and bold capital letters to represent matrices. Denote by \mathbf{u}_i the *i*-th element of vector \mathbf{u} . For two vectors **u**, $\mathbf{v} \in \mathbb{R}^d$, we say $\mathbf{u} \leq \mathbf{v}$ iff $\mathbf{u}_i \leq \mathbf{v}_i$ for $\forall i \in [d]$. Denote by $\mathbf{v} \cdot \mathbf{u}$ their inner product, and by proj_ $\mathbf{u}(\mathbf{v}) = \frac{\mathbf{v} \cdot \mathbf{u}}{\mathbf{u} \cdot \mathbf{u}} \mathbf{u}$ the projection of \mathbf{v} onto the line spanned by \mathbf{u} .

126 3.2 PRELIMINARIES AND DEFINITIONS

Since our methodology is based on techniques in multi-objective minimization (MoM), we first review some concepts from MoM, particularly the multiple gradient descent algorithm (MGDA).

130 3.2.1 Multi-Objective Minimization for Fairness

Denote by $f(\theta) = \{f_k(\theta)\}_{k \in [K]}$ the set of local clients' loss functions; the aim of MoM is to solve

$$\boldsymbol{\theta}^* = \arg\min_{\boldsymbol{\theta}} \boldsymbol{f}(\boldsymbol{\theta}), \tag{1}$$

where the minimization is performed w.r.t. the *partial ordering*. Finding θ^* could enforce fairness among the users since by setting setting $\theta = \theta^*$, it is not possible to reduce any of the local objective functions f_k without increasing at least another one. Here, θ^* is called a Pareto-optimal solution of Equation (1). Although finding Pareto-optimal solutions can be challenging, there are several methods to identify the Pareto-stationary solutions instead, which are defined as follows:

141 **Definition 3.1.** Pareto-stationary Mukai (1980): The vector θ^* is said to be Pareto-stationary iff 142 there exists a convex combination of the gradient-vectors $\{g_k(\theta^*)\}_{k\in[K]}$ which is equal to zero; that 143 is, $\sum_{k=1}^{K} \lambda_k g_k(\theta^*) = 0$, where $\lambda \ge 0$, and $\sum_{k=1}^{K} \lambda_k = 1$.

Lemma 3.2. Mukai (1980) Any Pareto-optimal solution is Pareto-stationary. On the other hand, if all $\{f_k(\theta)\}_{k \in [K]}$'s are convex, then any Pareto-stationary solution is weakly Pareto optimal².

There are many methods in the literature to find Pareto-stationary solutions among which MGDA is a popular one Mukai (1980); Fliege & Svaiter (2000); Désidéri (2012).

149 MGDA adaptively tunes $\{\lambda_k\}_{k \in [K]}$ by finding the minimal-norm element of the convex hull of the 150 gradient vectors defined as follows (we drop the dependence of \mathbf{g}_k to $\boldsymbol{\theta}^t$ for ease of notation hereafter)

$$\mathcal{G} = \{ \boldsymbol{g} \in \mathbb{R}^d | \boldsymbol{g} = \sum_{k=1}^K \lambda_k \boldsymbol{g}_k; \ \lambda_k \ge 0; \ \sum_{k=1}^K \lambda_k = 1 \}.$$
(2)

Denote the minimal-norm element of \mathcal{G} by $\mathfrak{d}(\mathcal{G})$. Then, either (i) $\mathfrak{d}(\mathcal{G}) = 0$, and therefore based on Lemma 3.2 $\mathfrak{d}(\mathcal{G})$ is a Pareto-stationary point; or (ii) $\mathfrak{d}(\mathcal{G}) \neq 0$ and the direction of $-\mathfrak{d}(\mathcal{G})$ is a common descent direction for all the objective functions $\{f_k(\theta)\}_{k \in [K]}$ Désidéri (2009), meaning that all the directional derivatives $\{\mathfrak{g}_k \cdot \mathfrak{d}(\mathcal{G})\}_{k \in [K]}$ are positive. Having positive directional derivatives is a *necessary* condition to ensure that the common direction is descent for all the objective functions.

¹⁶⁰ 161

 $^{{}^{2}\}theta^{*}$ is called a weakly Pareto-optimal solution of Equation (1) if there does not exist any θ such that $f(\theta) < f(\theta^{*})$; meaning that, it is not possible to improve *all* of the objective functions in $f(\theta^{*})$. Obviously, any Pareto optimal solution is also weakly Pareto-optimal but the converse may not hold.

162 3.2.2 NEWTON-TYPE METHODS

164 First-order FL methods face challenges with slow convergence, measured in terms of the number of iterations or communication rounds. Additionally, their convergence speed is intricately linked to 165 the condition number, influenced by factors such as the model's structure, choice of loss function, 166 and distribution of training data. In contrast, second-order methods exhibit significantly faster 167 performance due to their additional computational effort in estimating the local curvature of the loss 168 landscape. This, in turn, yields faster and more adaptive update directions. Despite requiring more 169 computations per communication round, second-order methods achieve fewer communication rounds. 170 In the context of FL, where communication often poses a bottleneck rather than computation, the 171 appeal of second-order methods has grown. Notably, the Newton's direction is obtained as 172

173 174

175 176

177

178

186

187

$\mathfrak{d}_N = -(\nabla^2 f(\boldsymbol{\theta}))^{-1} \nabla f(\boldsymbol{\theta}). \tag{3}$

4 MOTIVATION AND METHODOLOGY

We discuss our motivation in Section 4.1 based on which we elaborate on the inner-working of DQN-FL in Section 4.2.

179 180 4.1 Motivation

We begin with finding out how much the local loss function $f_k(\cdot)$, $k \in [K]$, changes when the server updates the global model as $\theta^{t+1} = \theta^t - \eta^t \mathfrak{d}^t$ at round t. In other words, we want to determine the rate of change $\Delta f_k(\theta^t) \triangleq f_k(\theta^{t+1}) - f_k(\theta^t)$ for the local loss functions. To do this, by writing the first-order Taylor expansion for the local loss function $f_k(\cdot)$, we obtain:

$$f_k(\boldsymbol{\theta}^{t+1}) = f_k(\boldsymbol{\theta}^t - \eta^t \boldsymbol{\mathfrak{d}}^t) \approx f_k(\boldsymbol{\theta}^t) - \eta^t \boldsymbol{\mathfrak{g}}_k^t \cdot \boldsymbol{\mathfrak{d}}^t$$
(4)

$$\Rightarrow \quad \Delta f_k(\boldsymbol{\theta}^t) \approx -\eta^t \boldsymbol{\mathfrak{g}}_k^t \cdot \boldsymbol{\mathfrak{d}}^t. \tag{5}$$

As per Equation (5), $f_k(\cdot)$ changes by amount of $-\eta^t \mathbf{g}_k^t \cdot \mathbf{d}^t$ when the server updates the global model. Hence, if $\mathbf{g}_k^t \cdot \mathbf{d} \ge 0$, the global updating direction is descent for client k, and $\Delta f_k(\boldsymbol{\theta}^t) \le 0$.

Nevertheless, updating toward a descent direction does not guarantee any meaningful convergence.
 Indeed, what can guarantee the convergence of GD-like algorithms is the rate of change in the loss function in each iteration³. This is in fact what makes the second-order methods to converge faster as the rate of change in the loss functions is automatically determined by the Hessian matrix.

This motivates us to see how the server can update the global model such that the rate of change in the local loss functions is the same as that when local clients update their local loss function using second-order methods.

Specifically, let d_k^t denote the rate of change in local loss function f_k when it updates its local model using Newton method; then, we have

$$d_k^t = \mathbf{g}_k^t \cdot \mathbf{\mathfrak{d}}_N = \mathbf{g}_k^t \cdot \left((\mathbf{H}_k^t)^{-1} \mathbf{\mathfrak{g}}_k^t \right) = \left(\mathbf{\mathfrak{g}}_k^t \right)^T (\mathbf{H}_k^t)^{-1} \mathbf{\mathfrak{g}}_k^t.$$
(6)

Our goal is to assist the server in updating the global model such that, after the update, the rate of change for client k becomes d_k^t . Achieving this is not a straightforward task. In the following section, we derive a closed-form solution to meet this criterion.

4.2 Methodology

Our method is partially inspired from MGDA algorithm, but incorporates several key modifications.
 Specifically, our approach comprises two stages: (i) gradient orthogonalization with a tailored scaling strategy; and (ii) finding the optimal weights to combine these orthogonal gradients.

210 211

205

206

4.2.1 Stage 1, GRADIENT ORTHOGONALIZATION

The clients send the local gradients $\{\mathbf{g}_k\}_{k \in [K]}$ to the server, and then the server first generates a mutually orthogonal ⁴ set $\{\tilde{\mathbf{g}}_k\}_{k \in [K]}$ that spans the same *K*-dimensional subspace in \mathbb{R}^d as that

³If $f_k(\cdot)$ is L-smooth, the convergence of gradient descent algorithm is guaranteed for $\eta^t \in [0, \frac{2}{L}]$.

⁴Here, orthogonality is in the sense of standard inner product in Euclidean space.

spanned by $\{g_k\}_{k \in [K]}$. To this aim, the server exploits a modified Gram–Schmidt orthogonalization process over $\{g_k\}_{k \in [K]}$ in the following manner ⁵

$$\tilde{\mathbf{g}}_1 = \mathbf{g}_1/d_1^t,\tag{7}$$

219 220 221

222

228 229

230

231

241

242

253 254 255

256 257

266

267 268

$$\tilde{\mathbf{g}}_{k} = \frac{\mathbf{g}_{k} - \sum_{i=1}^{k-1} \operatorname{proj}_{\tilde{\mathbf{g}}_{i}}(\mathbf{g}_{k})}{d_{k}^{t} - \sum_{i=1}^{k-1} \frac{\mathbf{g}_{k} \cdot \tilde{\mathbf{g}}_{i}}{\tilde{\mathbf{g}}_{i} \cdot \tilde{\mathbf{g}}_{i}}}, \text{ for } k = 2, \dots, K,$$
(8)

where $\gamma > 0$ is a scalar. Note that the orthogonalization approach in *stage* 1 is feasible if we assume that the *K* gradient vectors $\{\mathbf{g}_k\}_{k \in [K]}$ are linearly independent. Indeed, this assumption is reasonable considering that (i) the gradient vectors $\{\mathbf{g}_k\}_{k \in [K]}$ are *K* vectors in *d*-dimensional space, and d >> K for the DNNs⁶; and that (ii) the random nature of the gradient vectors due to the non-iid distributions of the local datasets.

4.2.2 *Stage* 2, FINDING OPTIMAL WEIGHTS

In this stage, we aim to find the minimum-norm vector in the convex hull of the *orthogonal* gradients found in *Stage (I)*. First, denote by $\tilde{\mathcal{G}}$ the convex hull of gradient vectors $\{\tilde{\mathbf{g}}_k\}_{k \in [K]}$; that is,

$$ilde{\mathcal{G}} = \{ oldsymbol{g} \in \mathbb{R}^d | oldsymbol{g} = \sum_{k=1}^K \lambda_k ilde{oldsymbol{g}}_k; \; \lambda_k \geq 0; \; \sum_{k=1}^K \lambda_k = 1 \}.$$

In the following, we find the minimal-norm element in $\tilde{\mathcal{G}}$, and then we show that this element is a descent direction for all the objective functions.

²³⁸ Denote by λ^* the weights corresponding to the minimal-norm vector in $\tilde{\mathcal{G}}$. To find the weight vector λ^* , we solve

$$\boldsymbol{g}^* = \arg\min_{\boldsymbol{g}\in\mathcal{G}} \|\boldsymbol{g}\|_2^2, \tag{9}$$

243 which accordingly finds λ^* . For an element $g \in \mathcal{G}$, we have

$$\|\boldsymbol{g}\|_{2}^{2} = \|\sum_{k=1}^{K} \lambda_{k} \tilde{\boldsymbol{\mathfrak{g}}}_{k}\|_{2}^{2} = \sum_{k=1}^{K} \lambda_{k}^{2} \|\tilde{\boldsymbol{\mathfrak{g}}}_{k}\|_{2}^{2},$$
(10)

where we used the fact that $\{\tilde{\mathbf{g}}_k\}_{k \in [K]}$ are orthogonal.

To solve Equation (9), we first ignore the inequality $\lambda_k \ge 0$, for $k \in [K]$, and then we observe that it is automatically satisfied. Thus, we make the following Lagrangian to solve the minimization problem in Equation (9):

Hence, $\frac{\partial L}{\partial \lambda_k} = 2\lambda_k \|\tilde{\mathbf{g}}_k\|_2^2 - \alpha$; and by setting this equation to zero we obtain

$$\lambda_k^* = \frac{\alpha}{2\|\tilde{\mathbf{g}}_k\|_2^2}.$$
(11)

On the other hand, since $\sum_{k=1}^{K} \lambda_k = 1$, from Equation (11) we have $\alpha = \frac{2}{\sum_{k=1}^{K} \frac{1}{\|\tilde{\mathfrak{g}}_k\|_2^2}}$ from which the optimal λ^* is obtained as follows

$$\lambda_k^* = \frac{1}{\|\tilde{\mathbf{g}}_k\|_2^2 \sum_{k=1}^K \frac{1}{\|\tilde{\mathbf{g}}_k\|_2^2}}, \quad \text{for } k \in [K].$$
(12)

Note that $\lambda_k^* > 0$, and therefore the minimum norm vector we found belongs to \mathcal{G} . Using the λ^* found in equation 12, we can calculate $\mathfrak{d}^t = \sum_{k=1}^K \lambda_k^* \tilde{\mathfrak{g}}_k$ as the minimum norm element in the convex hull $\tilde{\mathcal{G}}$.

Theorem 4.1. If the server updates the model toward $\mathfrak{d}^t = \sum_{k=1}^K \lambda_k^* \tilde{\mathfrak{g}}_k$, the rate of change for client k is proportional to d_k^t , $\forall k \in [K]$.

⁵The reason for such normalization will be clarified later.

⁶Also, note that to tackle non-iid distribution of user-specific data, it is a common practice that server selects a different subset of clients in each round McMahan et al. (2017).

 Proof. We shall find the directional derivative of loss function $f_k, \forall k \in [K]$, over \mathfrak{d}^t :

$$\mathbf{\mathfrak{g}}_{k} \cdot \mathbf{\mathfrak{d}}^{t} = \left(\tilde{\mathbf{\mathfrak{g}}}_{k} \left(d_{k}^{t} - \sum_{i=1}^{k-1} \frac{\mathbf{\mathfrak{g}}_{k} \cdot \tilde{\mathbf{\mathfrak{g}}}_{i}}{\tilde{\mathbf{\mathfrak{g}}}_{i} \cdot \tilde{\mathbf{\mathfrak{g}}}_{i}} \right) + \sum_{i=1}^{k-1} \operatorname{proj}_{\tilde{\mathbf{\mathfrak{g}}}_{i}}(\mathbf{\mathfrak{g}}_{k}) \right) \cdot \left(\sum_{i=1}^{K} \lambda_{i}^{*} \tilde{\mathbf{\mathfrak{g}}}_{i}\right)$$
(13)

$$=\lambda_{k}^{*}\|\tilde{\mathbf{g}}_{k}\|_{2}^{2}\left(d_{k}^{t}-\sum_{i=1}^{k-1}\frac{\mathbf{g}_{k}\cdot\tilde{\mathbf{g}}_{i}}{\tilde{\mathbf{g}}_{i}\cdot\tilde{\mathbf{g}}_{i}}\right)+\sum_{i=1}^{k-1}\frac{\mathbf{g}_{k}\cdot\tilde{\mathbf{g}}_{i}}{\tilde{\mathbf{g}}_{i}\cdot\tilde{\mathbf{g}}_{i}}\lambda_{i}^{*}\|\tilde{\mathbf{g}}_{i}\|_{2}^{2}$$
(14)

$$= \frac{\alpha}{2} \left(d_k^t - \sum_{i=1}^{k-1} \frac{\mathbf{g}_k \cdot \tilde{\mathbf{g}}_i}{\tilde{\mathbf{g}}_i \cdot \tilde{\mathbf{g}}_i} \right) + \frac{\alpha}{2} \sum_{i=1}^{k-1} \frac{\mathbf{g}_k \cdot \tilde{\mathbf{g}}_i}{\tilde{\mathbf{g}}_i \cdot \tilde{\mathbf{g}}_i}$$
(15)

$$= \frac{\alpha}{2} d_k^t = \frac{d_k^t}{\sum_{k=1}^K \frac{1}{\|\hat{\mathbf{g}}_k\|_2^2}} > 0, \tag{16}$$

where (i) Equation (13) is obtained by using definition of $\tilde{\mathbf{g}}_k$ in Equation (8), (ii) Equation (14) follows from the orthogonality of $\{\tilde{\mathbf{g}}_k\}_{k=1}^K$ vectors, and (iii) Equation (15) is obtained by using Equation (11).

Hence, to realize a rate of change similar to the Newton step, at iteration t, the server set the global learning rate as $\eta = \sum_{k=1}^{K} \frac{1}{\|\tilde{\mathbf{g}}_{k}\|_{2}^{2}}$, and update the global model as:

$$\boldsymbol{\theta}^{t+1} = \boldsymbol{\theta}^t - \eta^t \boldsymbol{\mathfrak{d}}^t = \boldsymbol{\theta}^t - \sum_{k=1}^K \frac{1}{\|\tilde{\boldsymbol{\mathfrak{g}}}_k\|_2^2} \boldsymbol{\mathfrak{d}}^t.$$
(17)

To summarize, updating the global model as in Equation (17) provides two key advantages: (i) All local loss decreases (as shown by the inequality in Equation (16)); (ii) The rate of change for each local loss function aligns with that of the Newton method.

Similarly to the conventional GD, we note that updating the global model as equation 17 is a *necessary* condition to have $f(\theta^{t+1}) \leq f(\theta^t)$. In Theorem 4.2 whose proof is differed to Appendix A, we state the *sufficient* condition to satisfy $f(\theta^{t+1}) \leq f(\theta^t)$.

Theorem 4.2. Assume that $f = {f_k}_{k \in [K]}$ are L-Lipschitz smooth. If the step-size $\eta^t =$ $\sum_{k=1}^{K} \frac{1}{\|\tilde{\mathbf{g}}_{k}\|_{2}^{2}} \in [0, \frac{2}{L} \min\{d_{k}^{t}\}_{k \in [K]}], \text{ then } \boldsymbol{f}(\boldsymbol{\theta}^{t+1}) \leq \boldsymbol{f}(\boldsymbol{\theta}^{t}), \text{ and equality is achieved iff } \boldsymbol{\mathfrak{d}}^{t} = \boldsymbol{0}.$

4.3 DQN-FED ALGORITHM

Since Newton's method requires the computation of the inverse Hessian matrix, which is computationally expensive, we employ quasi-Newton methods that approximate the inverse of the Hessian using gradient information. The BFGS algorithm Broyden (1970) is one such approach. Let B_k^t denote the matrix obtained using BFGS algorithm, where $\hat{B}_k^t \approx (H_k^t)^{-1}$. Using B_k^t , d_k^t in Equation (6) can be approximated by

$$\tilde{d}_k^t = (\mathbf{g}_k^t)^T \boldsymbol{B}_k^t \mathbf{g}_k^t.$$
(18)

Lastly, similar to many recent FL algorithms McMahan et al. (2017); Li et al. (2019a), we allow each client to perform a couple of local epochs e. We summarize DQN-Fed in Algorithm 1.

4.4 **CONVERGENCE RESULTS**

In the following, we prove the convergence guarantee of DQN-Fed based on how the clients update the local models: (i) using SGD with e = 1, (ii) using GD with e > 1, and (iii) using GD with e = 1. Naturally, the strongest convergence guarantee is provided for the third scenario.

Theorem 4.3 (e = 1 & local SGD). Assume that $f = {f_k}_{k \in [K]}$ are l-Lipschitz continuous and L-Lipschitz smooth, and that the global step-size η^t satisfies the following three conditions:

Input:	Number of global epochs T, global learning rate η^t , number of local epochs	E, local
dataset	ts $\{\mathcal{D}_k\}_{k\in K}$.	
for $t =$	$=0,1,\ldots,T-1$ do	
Sei	rver randomly selects a subset of devices S^t and sends θ^t to them.	
for	$device k \in S^t$ in parallel do	
	Set $\hat{\theta}_k^0 = \theta^t$ and $\hat{\theta}_k^{-1} = \theta^{t-1}$	
	for $e = 0, 1,, E - 1$ do	
	Perform BFGS algorithm as follows	
	Set $\mathbf{s}_k^e = \hat{\boldsymbol{\theta}}_k^e - \hat{\boldsymbol{\theta}}_k^{e-1}$, and $\mathbf{y}_k^e = \nabla f(\hat{\boldsymbol{\theta}}_k^e) - \nabla f(\hat{\boldsymbol{\theta}}_k^{e-1})$.	
	Iteratively update matrix \mathbf{B}_k^{e+1} using information from $\mathbf{B}_k^e, \mathbf{s}_k^e, \mathbf{y}_k^e$ accord	ding to:
		-
	$\mathbf{B}_k^{e+1} = \mathbf{B}_k^e - rac{\mathbf{B}_k^e \mathbf{s}_k^e (\mathbf{s}_k^e)^\intercal \mathbf{B}_k^e}{(\mathbf{s}_k^t)^\intercal \mathbf{B}_k^e \mathbf{s}_k^e} + rac{\mathbf{y}_k^e (\mathbf{y}_k^e)^\intercal}{(\mathbf{s}_k^e)^\intercal \mathbf{y}_k^e}.$	(19)
	end	
	Use \mathbf{B}_k^E to calculate \tilde{d}_k^t from Equation (18).	
	Send local gradient $\mathbf{g}_k = \nabla f(\hat{\boldsymbol{\theta}}_k^e)$ and \tilde{d}_k^t to the server.	
en	\mathbf{d}	
Sei	rver finds $\{\tilde{\mathfrak{g}}_k\}_{k\in[K]}$ form Equations (7) and (8).	
Sei	rver finds λ^* from Equation (12).	
Se	rver calculates $\boldsymbol{\vartheta}^t := \sum_{i=1}^{K} \lambda_i^* \tilde{\mathfrak{q}_i}$.	
Set	rver updates the global model as $\theta^{t+1} \leftarrow \theta^t - n^t \mathfrak{d}^t$	
end		
-		

(i) $\eta^t \in (0, \frac{1}{2L}]$, (ii) $\lim_{T\to\infty} \sum_{t=0}^T \eta^t \to \infty$, and (iii) $\lim_{T\to\infty} \sum_{t=0}^T \eta^t \sigma_t < \infty$; where $\sigma_t^2 = \mathbf{E}[\|\tilde{\mathbf{g}}\boldsymbol{\lambda}^* - \tilde{\mathbf{g}}_s \boldsymbol{\lambda}_s^*\|]^2$ is the variance of stochastic common descent direction. Then

$$\lim_{T \to \infty} \min_{t=0,\dots,T} \mathbf{E}[\|\boldsymbol{\mathfrak{d}}^t\|] \to 0.$$
⁽²⁰⁾

Theorem 4.4 (e > 1 & local GD). Assume that $f = {f_k}_{k \in [K]}$ are l-Lipschitz continuous and L-Lipschitz smooth. Denote by η^t and η the global and local learning rates, respectively. Also, define $\zeta^t = \| \lambda^* - \lambda_e^* \|$, where λ_e^* is the optimum weights obtained from pseudo-gradients after e local epochs. We have

$$\lim_{T \to \infty} \min_{t=0,\dots,T} \|\boldsymbol{\mathfrak{d}}^t\| \to 0, \tag{21}$$

if the following conditions are satisfied: (i) $\eta^t \in (0, \frac{1}{2L}]$, (ii) $\lim_{T\to\infty} \sum_{t=0}^T \eta^t \to \infty$, (iii) $\lim_{t\to\infty} \eta^t \to 0$, (iv) $\lim_{t\to\infty} \eta \to 0$, and (v) $\lim_{t\to\infty} \zeta^t \to 0$.

Before introducing Theorem 4.5, we first introduce some notations. Denote by ϑ the Pareto-stationary solution set⁷ of minimization problem $\arg \min_{\theta} f(\theta)$. Then, denote by θ^* the projection of θ^t onto the set ϑ ; that is, $\theta^* = \arg \min_{\theta \in \vartheta} \|\theta^t - \theta\|_2^2$.

Theorem 4.5 (e = 1 & local GD). Assume that $f = {f_k}_{k \in [K]}$ are l-Lipschitz continuous and σ -convex, and that the global step-size η^t satisfies the following two conditions: (i) $\lim_{t\to\infty} \sum_{j=0}^t \eta_j \to t$ ∞ , and (ii) $\lim_{t\to\infty} \sum_{i=0}^t \eta_i^2 < \infty$. Then almost surely $\theta^t \to \theta^*$; that is,

$$\mathbb{P}\left(\lim_{t\to\infty}\left(\boldsymbol{\theta}^t-\boldsymbol{\theta}^*\right)=0\right)=1,\tag{22}$$

where $\mathbb{P}(E)$ denotes the probability of event E.

⁷In general, the Pareto-stationary solution of multi-objective minimization problem forms a set with cardinality of infinity Mukai (1980).

The proofs for Theorems 4.3 to 4.5 are provided in Appendices B.1 to B.3, respectively. Note that all the Theorems 4.3 to 4.5 provide some types of convergence to a Pareto-optimal solution of optimization problem in Equation (1). Specifically, diminishing $\boldsymbol{\vartheta}^t$ in Theorems 4.3 and 4.4 implies that we are reaching to a Pareto-optimal point Désidéri (2009). On the other hand, Theorem 4.5 explicitly provides this convergence guarantee in an almost surely fashion.

In addition, the following theorem shows that DQN-Fed has a *linear-quadratic* convergence rate (for the full-version of the theorem and its proof refer to Appendix D).

Theorem 4.6 (Convergence rate of DQN-Fed). Assume that the global loss function is twice continuously differentiable, *L*-Lipschitz gradient (*L*-smooth) and λ -strongly convex. In addition, assume that the matrix \mathbf{B}_t^{-1} is a δ -approximate of true inverse Hessian \mathbf{H}_t^{-1} ; that is $\|\mathbf{B}_t^{-1} - \mathbf{H}_t^{-1}\| \leq \delta \|\mathbf{H}_t^{-1}\|$. Then, the convergence is either *quadratic*, with $T_{\epsilon} = \mathcal{O}\left(\log \log \frac{1}{\epsilon}\right)$, or *linear*, with $T_{\epsilon} = \mathcal{O}\left(\frac{1}{\log(\frac{1}{\epsilon \lambda})} \log \frac{1}{\epsilon}\right)$, where λ , *L* and δ are constants.

386

387

388

389 390

391 392 393

394 395

397

398

5 EXPERIMENTS

In this section, we conclude the paper by presenting a series of experiments to demonstrate the performance of DQN-Fed. We also conduct a comparative analysis to assess its effectiveness against state-of-the-art alternatives using various performance metrics.

399 • Datasets: We conduct a comprehensive set of experiments across seven datasets. In this section, 400 we present results for four datasets: CIFAR-{10, 100} Krizhevsky et al. (2009), FEMNIST Caldas 401 et al. (2018), and Shakespeare McMahan et al. (2017). Results for Fashion MNIST Xiao et al. (2017), 402 TinyImageNet Le & Yang (2015), and CINIC-10 Darlow et al. (2018) are discussed in Appendix I. To demonstrate DON-Fed's effectiveness across different FL scenarios, we examine two FL setups for 403 each dataset in this section. We provide some experimental analysis in Appendix H where we show 404 that DQN-Fed converges faster than the first-order FL methods. Furthermore, we evaluate DQN-Fed's 405 performance on a real-world noisy dataset, Clothing1M Xiao et al. (2015), in Appendix K. 406

Benchmarks: We compare the performance of DQN-Fed against some fair first-order FL and some second-order FL methods. The fair FL algorithms include q-FFL Li et al. (2019a), TERM Li et al. (2020), FedMGDA+ Hu et al. (2022), Ditto Li et al. (2021), FedLF Pan et al. (2024), FedHEAL Chen et al. (2024), and conventional FedAvg McMahan et al. (2017); and also second-order FL methods include FedNL (Safaryan et al., 2022) and FedNew (Elgabli et al., 2022).

It is worth noting that we conduct a grid-search to find the best hyper-parameters for each of the
benchmark methods including DQN-Fed in our experiments. The details of this hyper-parameter
tuning are reported in Appendix J.

• Performance metrics: Denote by a_k the prediction accuracy on device k. We use $\bar{a} = \frac{1}{K} \sum_{k=1}^{K} a_k$ as the average test accuracy of the underlying FL algorithm, and use $\sigma_a = \sqrt{\frac{1}{K} \sum_{k=1}^{K} (a_k - \bar{a})^2}$ as the standard deviation of the accuracy across the clients. Furthermore, we report Worst 10% (5%) and Best 10% (5%) accuracies as a common metric in fair FL algorithms Li et al. (2020).

• Notations: We use **bold** and <u>underlined</u> numbers to denote the best and second best performance, respectively. We use *e* and *K* to represent the number of local epochs and that of clients, respectively.

422 423

420

421

424 5.1 CIFAR-10

425 426 CIFAR-10 dataset Krizhevsky et al. (2009) has 50K training and 10K test images of size 32×32 427 labeled for 10 classes. The batch size is equal to 64 for both of the following setups.

• Setup 1: Following Wang et al. (2021b), we sort the dataset based on their classes, and then split them into 200 shards. Each client randomly selects two shards without replacement so that each has the same local dataset size. We use a feedforward neural network with 2 hidden layers. We fix e = 1and K = 100. We carry out 2000 rounds of communication, and sample 10% of the clients in each round. We run SGD on local datasets with stepsize $\eta = 0.1$.

• Setup 2: We distribute the dataset among the clients deploying Dirichlet allocation Wang et al. (2020) with $\beta = 0.5$. We use ResNet-18 He et al. (2016) with Group Normalization Wu & He (2018). We perform 100 communication rounds in each of which all clients participate. We set e = 1, K = 10 and $\eta = 0.01$.

Table 1: Test accuracy	v on CIFAR-10.	The reported results at	e averaged over 5 seeds.
fuolo f. fost accurac	y on on rine 10.	The reported results u	e averagea over 5 secas.

		Setup 1			Setup 2				
	Algorithm	ā	σ_a	W(5%)	B(5%)	$ \bar{a}$	σ_a	W(10%)	B(10%)
Naive First-Order	FedAvg	46.85	3.54	19.84	69.28	63.55	5.44	53.40	72.24
Fair First-Order	q-FFL FedMGDA FedHEAL TERM Ditto	46.30 45.34 46.40 47.11 46.31	$\begin{array}{c} \underline{3.27} \\ 3.37 \\ 3.61 \\ 3.66 \\ 3.44 \end{array}$	23.39 24.00 19.33 <u>28.21</u> 27.14	68.02 68.51 69.30 69.51 68.44	57.27 62.05 63.05 64.15 63.49	5.60 4.88 4.95 5.90 5.70	47.29 52.69 48.69 <u>56.21</u> 55.99	66.92 70.77 70.88 72.20 71.34
Second-Order	FedNL FedNew DQN-Fed	47.33 47.51 47.72	3.92 3.68 3.20	24.41 25.77 29.34	<u>69.52</u> 69.74 69.37	64.72 64.58 64.88	6.02 6.11 <u>4.90</u>	56.20 56.96 58.01	72.33 72.12 72.88

5.2 CIFAR-100

CIFAR-100 Krizhevsky et al. (2009) has the same number of samples as CIFAR-10, but comprises 100 classes compared to the 10 classes found in CIFAR-10.

The model for both setups is ResNet-18 He et al. (2016) with Group Normalization Wu & He (2018), where all clients participate in each round. We also set e = 1 and $\eta = 0.01$. The batch size is equal to 64. The results are reported in Table 2 for both of the following setups:

• Setup 1: We set K = 10 and $\beta = 0.5$ for Dirichlet allocation, and use 400 communication rounds.

• Setup 2: We set K = 50 and $\beta = 0.05$ for Dirichlet allocation, and use 200 communication rounds.

			Setup 1		Setup 2				
	Algorithm	ā	σ_a	W(10%)	B(10%)	ā	σ_a	W(10%)	B(10%)
Naive First-Order	FedAvg	30.05	4.03	25.20	40.31	20.15	6.40	11.20	33.80
Fair First-Order	q-FFL FedMGDA FedLF TERM Ditto	28.86 29.12 30.28 30.34 29.81	4.44 4.17 3.68 3.51 3.79	25.38 25.67 25.33 27.03 26.90	39.77 39.71 39.45 39.35 39.39	20.20 20.15 18.92 17.88 17.52	$ \begin{array}{r} 6.24 \\ 5.41 \\ \underline{4.90} \\ 5.98 \\ 5.65 \end{array} $	11.09 11.12 <u>11.29</u> 10.09 10.21	34.02 33.92 28.60 31.68 31.25
Fair First-Order	FedNL FedNew DQN-Fed	<u>31.58</u> 30.95 32.58	4.55 4.39 <u>3.60</u>	27.14 27.19 27.91	40.62 40.55 40.99	22.74 21.16 23.15	6.02 5.27 4.45	<u>12.15</u> 11.77 12.81	34.44 34.27 35.11

Table 2: Test accuracy on CIFAR-100. The reported results are averaged over 5 different seeds.

5.3 FEMNIST

FEMNIST (Federated Extended MNIST) Caldas et al. (2018) is a federated image dataset distributed over 3,550 devices which has 62 classes containing 28×28 -pixel images of digits (0-9) and English characters (A-Z, a-z). For implementation, we use a CNN model with 2 convolutional layers followed by 2 fully-connected layers. The batch size is 32, and e = 2 for both of the following setups:

• **FEMNIST-original:** We use the setting in Li et al. (2021), and randomly sample K = 500 devices and train models using the default data stored in each device.

• FEMNIST-skewed: K = 100. We sample 10 lower case characters ('a'-'j') from Extended MNIST (EMNIST), and randomly assign 5 classes to each of the 100 devices.

Consistent with Li et al. (2019a), we use two other fairness metrics for this dataset: (i) the angle between the accuracy distribution and the all-ones vector 1 denoted by Angle (°), and (ii) the KL divergence between the normalized accuracy a and uniform distribution u denoted by KL (a||u). Results for both setups are reported in Table 3.

		FEMNIST-original					FEMN	IST-skewe	ed
	Algorithm	ā	σ_a	Ang (°)	$\operatorname{KL}\left(a\ u ight)$	ā	σ_a	Ang (°)	KL $(a u)$
Naive First-Order	FedAvg	80.42	11.16	10.18	0.017	79.24	22.30	12.29	0.054
Fair First-Order	q-FFL FedMGDA TERM FedLF Ditto	80.91 81.00 81.08 82.45 83.77	10.62 10.41 10.32 <u>9.85</u> 10.13	9.71 10.04 9.15 <u>9.01</u> 9.34	0.016 0.016 0.015 <u>0.012</u> 0.014	84.65 85.41 84.29 85.21 92.51	18.56 17.36 13.88 14.92 14.32	12.01 11.63 11.27 11.44 11.45	0.038 0.032 0.025 0.027 <u>0.022</u>
Fair First-Order	FedNL FedNew DQN-Fed	84.21 <u>84.25</u> 85.15	11.22 10.88 9.58	10.07 9.78 8.14	0.015 0.014 0.010	92.94 92.25 93.80	16.45 15.21 <u>13.91</u>	12.56 11.92 <u>11.41</u>	0.045 0.037 0.011

Table 3: Test accuracy on FEMNIST. The reported results are averaged over 5 different seeds.

5.4 TEXT DATA

We use The Complete Works of William Shakespeare McMahan et al. (2017) as the dataset, and train an RNN whose input is 80-character sequence to predict the next character. We use e = 1, and let all the devices participate in each round. The results are reported in Table 4 for the following setups:

• Setup 1: Following McMahan et al. (2017), we subsample 31 speaking roles, and assign each role to a client (K = 31) to complete 500 communication rounds. We use a model with two LSTM layers Hochreiter & Schmidhuber (1997) and one densely-connected layer. The initial $\eta = 0.8$ with decay rate of 0.95.

• Setup 2: Among the 31 speaking roles, the 20 ones with more than 10000 samples are selected, and assigned to 20 clients. We use an LSTM followed by a fully-connected layer. $\eta = 2$, and the number of communication is 100.

Table 4: Test accuracy on Shakespeare. The reported results are averaged over 5 different seeds.

		Setup 1				Setup 2			
	Algorithm	ā	σ_a	W(10%)	B(10%)	ā	σ_a	W(10%)	B(10%)
Naive First-Order	FedAvg	53.21	9.25	51.01	58.41	50.48	1.24	48.20	52.10
Fair First-Order	q-FFL FedMGDA FedLF TERM	53.90 53.08 54.58 54.16	$\begin{array}{ c c c }\hline 7.52 \\ 8.14 \\ 8.44 \\ 8.21 \end{array}$	51.52 52.84 52.87 52.09	58.47 58.51 59.84 59.15	50.72 50.41 52.45 52.17	$\begin{array}{c c} \underline{1.07} \\ 1.09 \\ 1.23 \\ 1.11 \end{array}$	48.90 48.18 50.02 49.14	52.29 51.99 54.17 53.62
	Ditto	60.74	8.32	53.57	64.92	53.12	1.20	50.94	55.23
Fair First-Order	FedNL FedNew DQN-Fed	60.25 60.59 61.65	8.24 7.55 6.55	53.15 53.18 53.79	64.15 64.09 <u>64.86</u>	52.24 52.49 <u>52.89</u>	1.25 1.19 0.98	50.77 50.82 51.02	54.41 54.36 <u>54.48</u>

5.5 ANALYSIS OF RESULTS

Based on the insights gleaned from Tables 1 to 4, several noteworthy observations emerge:

(i) Naive second-order FL methods, namely FedNL and FedNew, tend to train unfair models, despite achieving high average accuracy across clients.

(ii) Compared to the benchmark models, DQN-Fed consistently trains models that demonstrate significantly higher levels of fairness across clients.

(*iii*) The average accuracy of the model learned by DQN-Fed is higher compared to both first-order and second-order FL methods.

CONCLUSION

This paper introduced distributed quasi-Newton federated learning (DQN-Fed), a novel approach designed to ensure fairness while harnessing the fast convergence properties of quasi-Newton methods in FL setting. DQN-Fed aids the server in updating the global model by ensuring (i) all local loss functions decrease, promoting fairness; and (ii) the rate of change in local loss functions matches that of the quasi-Newton method. We prove the convergence of DQN-Fed and establish its linear-quadratic convergence rate. Furthermore, we validate DQN-Fed's effectiveness across various federated datasets, demonstrating its superiority over state-of-the-art fair FL methods.

540 REFERENCES

546

552

562

566

569

570

542	Solon Barocas, Moritz Hardt, and Arvind Narayanan. Fairness in machine learning. Nips tutorial, 1:
543	2017, 2017.

- Roberto Battiti. First-and second-order methods for learning: between steepest descent and newton's method. *Neural computation*, 4(2):141–166, 1992.
- 547 Charles G Broyden. The convergence of a class of double-rank minimization algorithms: 2. the new algorithm. *IMA journal of applied mathematics*, 6(3):222–231, 1970.
- Sebastian Caldas, Sai Meher Karthik Duddu, Peter Wu, Tian Li, Jakub Konečnỳ, H Brendan McMa han, Virginia Smith, and Ameet Talwalkar. Leaf: A benchmark for federated settings. *arXiv preprint arXiv:1812.01097*, 2018.
- Yuhang Chen, Wenke Huang, and Mang Ye. Fair federated learning under domain skew with local consistency and domain diversity. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 12077–12086, 2024.
- Sen Cui, Weishen Pan, Jian Liang, Changshui Zhang, and Fei Wang. Addressing algorithmic disparity
 and performance inconsistency in federated learning. *Advances in Neural Information Processing Systems*, 34:26091–26102, 2021.
- ⁵⁵⁹ Nicolò Dal Fabbro, Subhrakanti Dey, Michele Rossi, and Luca Schenato. Shed: A newton-type algorithm for federated learning based on incremental hessian eigenvector sharing. *Automatica*, 160:111460, 2024.
- Luke N Darlow, Elliot J Crowley, Antreas Antoniou, and Amos J Storkey. Cinic-10 is not imagenet or cifar-10. *arXiv preprint arXiv:1810.03505*, 2018.
- ⁵⁶⁵ Jean-Antoine Désidéri. *Multiple-gradient descent algorithm (MGDA)*. PhD thesis, INRIA, 2009.
- Jean-Antoine Désidéri. Multiple-gradient descent algorithm (mgda) for multiobjective optimization.
 Comptes Rendus Mathematique, 350(5-6):313–318, 2012.
 - Wei Du, Depeng Xu, Xintao Wu, and Hanghang Tong. Fairness-aware agnostic federated learning. In Proceedings of the 2021 SIAM International Conference on Data Mining (SDM), pp. 181–189. SIAM, 2021.
- Hubert Eichner, Tomer Koren, Brendan McMahan, Nathan Srebro, and Kunal Talwar. Semi-cyclic stochastic gradient descent. In *International Conference on Machine Learning*, pp. 1764–1773. PMLR, 2019.
- Anis Elgabli, Chaouki Ben Issaid, Amrit Singh Bedi, Ketan Rajawat, Mehdi Bennis, and Vaneet
 Aggarwal. Fednew: A communication-efficient and privacy-preserving newton-type method for
 federated learning. In *International conference on machine learning*, pp. 5861–5877. PMLR, 2022.
- Jörg Fliege and Benar Fux Svaiter. Steepest descent methods for multicriteria optimization. *Mathematical methods of operations research*, 51(3):479–494, 2000.
- Avishek Ghosh, Raj Kumar Maity, and Arya Mazumdar. Distributed newton can communicate
 less and resist byzantine workers. *Advances in Neural Information Processing Systems*, 33: 18028–18038, 2020.
- Vipul Gupta, Avishek Ghosh, Michal Derezinski, Rajiv Khanna, Kannan Ramchandran, and Michael
 Mahoney. Localnewton: Reducing communication bottleneck for distributed learning. *arXiv preprint arXiv:2105.07320*, 2021.
- Shayan Mohajer Hamidi and EN-HUI YANG. Adafed: Fair federated learning via adaptive common descent direction. *Transactions on Machine Learning Research*, 2024. ISSN 2835-8856. URL https://openreview.net/forum?id=rFecyFpFUp.
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image
 recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*,
 pp. 770–778, 2016.

602

624

632

633

634

641

- Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. *Neural computation*, 9(8):
 1735–1780, 1997.
- Zeou Hu, Kiarash Shaloudegi, Guojun Zhang, and Yaoliang Yu. Federated learning meets multi objective optimization. *IEEE Transactions on Network Science and Engineering*, 2022.
- SHI Huaizhou, R Venkatesha Prasad, Ertan Onur, and IGMM Niemegeers. Fairness in wireless networks: Issues, measures and challenges. *IEEE Communications Surveys & Tutorials*, 16(1): 5–24, 2013.
- Tiansheng Huang, Weiwei Lin, Wentai Wu, Ligang He, Keqin Li, and Albert Y Zomaya. An
 efficiency-boosting client selection scheme for federated learning with fairness guarantee. *IEEE Transactions on Parallel and Distributed Systems*, 32(7):1552–1564, 2020a.
- Tiansheng Huang, Weiwei Lin, Li Shen, Keqin Li, and Albert Y Zomaya. Stochastic client selection
 for federated learning with volatile clients. *IEEE Internet of Things Journal*, 9(20):20055–20070,
 2022.
- Wei Huang, Tianrui Li, Dexian Wang, Shengdong Du, and Junbo Zhang. Fairness and accuracy in federated learning. *arXiv preprint arXiv:2012.10069*, 2020b.
- Jiawen Kang, Zehui Xiong, Dusit Niyato, Han Yu, Ying-Chang Liang, and Dong In Kim. Incentive
 design for efficient federated learning in mobile networks: A contract theory approach. In 2019
 IEEE VTS Asia Pacific Wireless Communications Symposium (APWCS), pp. 1–5. IEEE, 2019.
- Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images. 2009.
- Mounssif Krouka, Anis Elgabli, Chaouki Ben Issaid, and Mehdi Bennis. Communication-efficient
 federated learning: A second order newton-type method with analog over-the-air aggregation.
 IEEE Trans. Green Commun. Netw., 6(3):1862–1874, 2022.
- Tra Huong Thi Le, Nguyen H Tran, Yan Kyaw Tun, Minh NH Nguyen, Shashi Raj Pandey, Zhu Han, and Choong Seon Hong. An incentive mechanism for federated learning in wireless cellular networks: An auction approach. *IEEE Transactions on Wireless Communications*, 20(8):4874–4887, 2021.
- Ya Le and Xuan Yang. Tiny imagenet visual recognition challenge. *CS* 231N, 7(7):3, 2015.
- Yann LeCun, Léon Bottou, Yoshua Bengio, and Patrick Haffner. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998.
- Jian Li, Yong Liu, and Weiping Wang. Fedns: A fast sketching newton-type algorithm for federated
 learning. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, pp. 13509–
 13517, 2024.
 - Tian Li, Maziar Sanjabi, Ahmad Beirami, and Virginia Smith. Fair resource allocation in federated learning. In *International Conference on Learning Representations*, 2019a.
- Tian Li, Ahmad Beirami, Maziar Sanjabi, and Virginia Smith. Tilted empirical risk minimization. In
 International Conference on Learning Representations, 2020.
- Tian Li, Shengyuan Hu, Ahmad Beirami, and Virginia Smith. Ditto: Fair and robust federated learning through personalization. In *International Conference on Machine Learning*, pp. 6357–6368. PMLR, 2021.
 - Xiang Li, Kaixuan Huang, Wenhao Yang, Shusen Wang, and Zhihua Zhang. On the convergence of fedavg on non-iid data. *arXiv preprint arXiv:1907.02189*, 2019b.
- Lingjuan Lyu, Xinyi Xu, Qian Wang, and Han Yu. Collaborative fairness in federated learning.
 Federated Learning: Privacy and Incentive, pp. 189–204, 2020.
- Brendan McMahan, Eider Moore, Daniel Ramage, Seth Hampson, and Blaise Aguera y Arcas.
 Communication-efficient learning of deep networks from decentralized data. In *Artificial intelligence and statistics*, pp. 1273–1282. PMLR, 2017.

648 649 650	Quentin Mercier, Fabrice Poirion, and Jean-Antoine Désidéri. A stochastic multiple gradient descent algorithm. <i>European Journal of Operational Research</i> , 271(3):808–817, 2018.
651 652 653	Shayan Mohajer Hamidi and Oussama Damen. Fair wireless federated learning through the identification of a common descent direction. <i>IEEE Communications Letters</i> , 28(3):567–571, 2024. doi: 10.1109/LCOMM.2024.3350378.
654 655 656	Mehryar Mohri, Gary Sivek, and Ananda Theertha Suresh. Agnostic federated learning. In <i>International Conference on Machine Learning</i> , pp. 4615–4625. PMLR, 2019.
657 658	Hiroaki Mukai. Algorithms for multicriterion optimization. <i>IEEE transactions on automatic control</i> , 25(2):177–186, 1980.
659 660 661 662	Takayuki Nishio and Ryo Yonetani. Client selection for federated learning with heterogeneous resources in mobile edge. In <i>ICC 2019-2019 IEEE international conference on communications (ICC)</i> , pp. 1–7. IEEE, 2019.
663 664 665	Zibin Pan, Chi Li, Fangchen Yu, Shuyi Wang, Haijin Wang, Xiaoying Tang, and Junhua Zhao. Fedlf: Layer-wise fair federated learning. In <i>Proceedings of the AAAI Conference on Artificial Intelligence</i> , volume 38, pp. 14527–14535, 2024.
666 667 668	Boris Polyak and Andrey Tremba. New versions of newton method: step-size choice, convergence domain and under-determined equations. <i>Optimization Methods and Software</i> , 35(6):1272–1303, 2020.
669 670 671 672	Mher Safaryan, Rustem Islamov, Xun Qian, and Peter Richtarik. Fednl: Making newton-type methods applicable to federated learning. In <i>International Conference on Machine Learning</i> , pp. 18959–19010. PMLR, 2022.
673 674 675	Zhendong Song, Hongguang Sun, Howard H Yang, Xijun Wang, Yan Zhang, and Tony QS Quek. Reputation-based federated learning for secure wireless networks. <i>IEEE Internet of Things Journal</i> , 9(2):1212–1226, 2021.
676 677 678	Hongyi Wang, Mikhail Yurochkin, Yuekai Sun, Dimitris Papailiopoulos, and Yasaman Khazaeni. Federated learning with matched averaging. <i>arXiv preprint arXiv:2002.06440</i> , 2020.
679 680 681	Jianyu Wang, Zachary Charles, Zheng Xu, Gauri Joshi, H Brendan McMahan, Maruan Al-Shedivat, Galen Andrew, Salman Avestimehr, Katharine Daly, Deepesh Data, et al. A field guide to federated optimization. <i>arXiv preprint arXiv:2107.06917</i> , 2021a.
683 684 685	Shusen Wang, Fred Roosta, Peng Xu, and Michael W Mahoney. Giant: Globally improved approximate newton method for distributed optimization. In <i>Adv. Neural Inf. Process. Syst.</i> , pp. 2332–2342, 2018.
686 687	Zheng Wang, Xiaoliang Fan, Jianzhong Qi, Chenglu Wen, Cheng Wang, and Rongshan Yu. Federated learning with fair averaging. <i>arXiv preprint arXiv:2104.14937</i> , 2021b.
689 690	Yuxin Wu and Kaiming He. Group normalization. In <i>Proceedings of the European conference on computer vision (ECCV)</i> , pp. 3–19, 2018.
691 692 693	Han Xiao, Kashif Rasul, and Roland Vollgraf. Fashion-mnist: a novel image dataset for benchmarking machine learning algorithms. <i>arXiv preprint arXiv:1708.07747</i> , 2017.
694 695 696	Tong Xiao, Tian Xia, Yi Yang, Chang Huang, and Xiaogang Wang. Learning from massive noisy labeled data for image classification. In <i>Proceedings of the IEEE conference on computer vision and pattern recognition</i> , pp. 2691–2699, 2015.
697 698 699 700	Miao Yang, Ximin Wang, Hongbin Zhu, Haifeng Wang, and Hua Qian. Federated learning with class imbalance reduction. In 2021 29th European Signal Processing Conference (EUSIPCO), pp. 2174–2178. IEEE, 2021.
701	Dongdong Ye, Rong Yu, Miao Pan, and Zhu Han. Federated learning in vehicular edge computing: A selective model aggregation approach. <i>IEEE Access</i> , 8:23920–23935, 2020.

702 703 704	Jingfeng Zhang, Cheng Li, Antonio Robles-Kelly, and Mohan Kankanhalli. Hierarchically fair federated learning. <i>arXiv preprint arXiv:2004.10386</i> , 2020.
705	Jingwen Zhang, Yuezhou Wu, and Rong Pan. Incentive mechanism for horizontal federated learning based on reputation and reverse auction. In <i>Proceedings of the Web Conference</i> 2021, pp. 947–956
706	2021
707	
708	Yuchen Zhang and Xiao Lin. Disco: Distributed optimization for self-concordant empirical loss. In
709	<i>ICML</i> , pp. 362–370. PMLR, 2015.
710	
710	
712	
714	
715	
716	
717	
718	
719	
720	
721	
722	
723	
724	
725	
726	
727	
728	
729	
730	
731	
732	
733	
734	
735	
730	
730	
730	
740	
741	
742	
743	
744	
745	
746	
747	
748	
749	
750	
751	
752	
753	
754	
755	

A PROOF OF THEOREM 4.2

Proof. If all the $\{f_k\}_{k \in [K]}$ are L-smooth, then

$$\boldsymbol{f}(\boldsymbol{\theta}^{t+1}) \leq \boldsymbol{f}(\boldsymbol{\theta}^{t}) + \boldsymbol{\mathfrak{g}}^{T}(\boldsymbol{\theta}^{t+1} - \boldsymbol{\theta}^{t}) + \frac{L}{2} \|\boldsymbol{\theta}^{t+1} - \boldsymbol{\theta}^{t}\|_{2}^{2}.$$
(23)

Now, for client $k \in [K]$, by using the update rule Equation (17) in Equation (23) we obtain

$$f_k(\boldsymbol{\theta}^{t+1}) \le f_k(\boldsymbol{\theta}^t) - \eta^t \mathfrak{g}_k \cdot \mathfrak{d}^t + (\eta^t)^2 \frac{L}{2} \| \mathfrak{d}^t \|_2^2.$$
(24)

To impose $f_k(\boldsymbol{\theta}^{t+1}) \leq f_k(\boldsymbol{\theta}^t)$, we should have

$$\eta^{t} \mathfrak{g}_{k} \cdot \mathfrak{d}^{t} \ge (\eta^{t})^{2} \frac{L}{2} \| \mathfrak{d}^{t} \|_{2}^{2}$$

$$\tag{25}$$

$$\Leftrightarrow \quad \mathfrak{g}_k \cdot \mathfrak{d}^t \ge \frac{\eta^t L}{2} \sum_{k=1}^K \frac{\|\tilde{\mathfrak{g}}_k\|_2^2}{\|\tilde{\mathfrak{g}}_k\|_2^4 \left(\sum_{i=1}^K \frac{1}{\|\tilde{\mathfrak{g}}_i\|_2^2}\right)^2} \tag{26}$$

$$\Leftrightarrow \quad \frac{\tilde{d}_{k}^{t}}{\sum_{k=1}^{K} \frac{1}{\|\tilde{\mathfrak{g}}_{k}\|_{2}^{2}}} \geq \frac{\eta^{t}L}{2} \frac{1}{\left(\sum_{k=1}^{K} \frac{1}{\|\tilde{\mathfrak{g}}_{k}\|_{2}^{2}}\right)^{2}} \sum_{k=1}^{K} \frac{1}{\|\tilde{\mathfrak{g}}_{k}\|_{2}^{2}} \tag{27}$$

$$\Leftrightarrow \quad \eta^t \le \frac{2}{L} \tilde{d}_k^t. \tag{28}$$

B CONVERGENCE OF DQN-FED

In the following, we provide three theorems to analyse the convergence of DQN-Fed under different scenarios. Specifically, we consider three cases: (i) Theorem B.1 considers e = 1 and using SGD for local updates, (ii) Theorem B.2 considers an arbitrary value for e and using GD for local updates, and (iii) Theorem B.4 considers e = 1 and using GD for local updates.

B.1 CASE 1: e = 1 & local SGD

Notations: We use subscript $(\cdot)_s$ to indicate a stochastic value. Using this notation for the values we introduced in the paper, our notations used in the proof of Theorem B.1 are summarized in Table 5.

Table 5: Notations used in Theorem B.1 for e = 1 & local SGD.

Notation	Description
$\mathfrak{g}_{k,s}$	Stochastic gradient vector of client k.
\mathfrak{g}_s	Matrix of Stochastic gradient vectors $[\mathfrak{g}_{1,s},\ldots,\mathfrak{g}_{K,s}]$.
$ ilde{\mathfrak{g}}_{k,s}$	Stochastic gradient vector of client k after orthogonalization process.
$ ilde{\mathfrak{g}}_s$	Matrix of orthogonalized Stochastic gradient vectors $[\tilde{\mathfrak{g}}_{1,s}, \dots, \tilde{\mathfrak{g}}_{K,s}]$.
$\lambda_{k,s}^*$	Optimum weights obtained from Equation equation 12 using Stochastic gradients $\tilde{\mathfrak{g}}_s$.
\mathfrak{d}_s	Optimum direction obtained using Stochastic $\tilde{\mathfrak{g}}_s$; that is, $\mathfrak{d}_s = \sum_{k=1}^K \lambda_{k,s}^* \tilde{\mathfrak{g}}_{k,s}$.

Theorem B.1. Assume that $f = \{f_k\}_{k \in [K]}$ are l-Lipschitz continuous and L-Lipschitz smooth, and that the step-size η^t satisfies the following three conditions: (i) $\eta^t \in (0, \frac{1}{2L}]$, (ii) $\lim_{T\to\infty} \sum_{t=0}^T \eta^t \to \infty$ and (iii) $\lim_{T\to\infty} \sum_{t=0}^T \eta^t \sigma_t < \infty$; where $\sigma_t^2 = \mathbf{E}[\|\tilde{\mathbf{g}}\boldsymbol{\lambda}^* - \tilde{\mathbf{g}}_s\boldsymbol{\lambda}_s^*\|]^2$ is the variance of stochastic common descent direction. Then

$$\lim_{T \to \infty} \min_{t=0,\dots,T} \mathbf{E}[\|\boldsymbol{\mathfrak{d}}^t\|] \to 0.$$
⁽²⁹⁾

⁸¹⁰ ⁸¹¹ ⁸¹² *Proof.* Since orthogonal vectors $\{\tilde{\mathfrak{g}}_k\}_{k\in[K]}$ span the same *K*-dimensional space as that spanned by gradient vectors $\{\mathfrak{g}_k\}_{k\in[K]}$, then

$$\exists \{\lambda'_k\}_{k \in [K]} \text{ s.t. } \boldsymbol{\mathfrak{d}} = \sum_{k=1}^K \lambda^*_k \tilde{\mathfrak{g}}_k = \sum_{k=1}^K \lambda'_k \mathfrak{g}_k = \boldsymbol{\mathfrak{g}} \boldsymbol{\lambda}'.$$
(30)

Similarly, for the stochastic gradients we have

825

835

842

843 844

850

851

852

853

854

861

863

$$\exists \{\lambda'_{k,s}\}_{k \in [K]} \text{ s.t. } \boldsymbol{\mathfrak{d}}_s = \sum_{k=1}^K \lambda^*_{k,s} \tilde{\boldsymbol{\mathfrak{g}}}_{k,s} = \sum_{k=1}^K \lambda'_{k,s} \boldsymbol{\mathfrak{g}}_{k,s} = \boldsymbol{\mathfrak{g}}_s \boldsymbol{\lambda}'_s.$$
(31)

Befine $\Delta_t = \mathfrak{g}\lambda' - \mathfrak{g}_s\lambda'_s = \tilde{\mathfrak{g}}\lambda^* - \tilde{\mathfrak{g}}_s\lambda^*_s$, where the last equality is due to the definitions in Equations (30) and (31).

We can find an upper bound for $f(\theta^{t+1})$ as follows

$$\boldsymbol{f}(\boldsymbol{\theta}^{t+1}) = \boldsymbol{f}(\boldsymbol{\theta}^t - \eta^t \boldsymbol{\mathfrak{d}}_t)$$
(32)

$$= \boldsymbol{f}(\boldsymbol{\theta}^{t} - \eta^{t} \sum_{k=1}^{K} \lambda_{k,s}^{*} \tilde{\boldsymbol{\mathfrak{g}}}_{k,s})$$
(33)

$$= \boldsymbol{f}(\boldsymbol{\theta}^{t} - \eta^{t}\boldsymbol{\mathfrak{g}}_{s}\boldsymbol{\lambda}_{s}^{\prime})$$
(34)

$$\leq \boldsymbol{f}(\boldsymbol{\theta}^{t}) - \eta^{t} \boldsymbol{\mathfrak{g}}^{T} \boldsymbol{\mathfrak{g}}_{s}^{T} \boldsymbol{\lambda}_{s}' + \frac{L(\eta^{t})^{2}}{2} \|\boldsymbol{\mathfrak{g}}_{s}^{T} \boldsymbol{\lambda}_{s}'\|^{2}$$
(35)

$$\leq \boldsymbol{f}(\boldsymbol{\theta}^{t}) - \eta^{t} \boldsymbol{g}^{T} \boldsymbol{g}^{T} \boldsymbol{\lambda}' + L(\eta^{t})^{2} \|\boldsymbol{g}^{T} \boldsymbol{\lambda}'\|^{2} + \eta^{t} \boldsymbol{g}^{T} \Delta_{t} + L(\eta^{t})^{2} \|\Delta_{t}\|^{2}$$
(36)

$$\leq \boldsymbol{f}(\boldsymbol{\theta}^{t}) - \eta^{t}(1 - L\eta^{t}) \|\boldsymbol{g}^{T}\boldsymbol{\lambda}'\|^{2} + l\eta^{t}\|\Delta_{t}\| + L(\eta^{t})^{2}\|\Delta_{t}\|^{2},$$
(37)

where equation 33 uses stochastic gradients in the updating rule of DQN-Fed, equation 34 is obtained from the definition in equation 31, equation 35 holds following the quadratic bound for smooth functions $f = \{f_k\}_{k \in [K]}$, and lastly equation 37 holds considering the Lipschits continuity of $f = \{f_k\}_{k \in [K]}$.

Assuming $\eta^t \in (0, \frac{1}{2L}]$ and taking expectation from both sides, we obtain:

 $\min_{t=0,\dots,T} \mathbf{E}[\|\boldsymbol{\mathfrak{d}}_t\|] \le \frac{\boldsymbol{f}(\boldsymbol{\theta}_0) - \mathbf{E}[\boldsymbol{f}(\boldsymbol{\theta}^{t+1})] + \sum_{t=0}^T \eta^t (l\sigma_t + L\eta^t \sigma_t^2)}{\frac{1}{2} \sum_{t=0}^T \eta^t}.$ (38)

Using the assumptions (i) $\lim_{T\to\infty} \sum_{j=0}^{T} \eta^t \to \infty$, and (ii) $\lim_{T\to\infty} \sum_{t=0}^{T} \eta^t \sigma_t < \infty$, the theorem will be concluded. Note that *vanishing* $\boldsymbol{\mathfrak{d}}^t$ implies reaching to a Pareto-stationary point of original MoM problem. Yet, the convergence rate is different in different scenarios as we see in the following theorems.

B.1.1 DISCUSSING THE ASSUMPTIONS

• The assumptions over the local loss functions: The two assumptions l-Lipschitz continuous and L-Lipschitz smooth over the local loss functions are two standard assumptions in FL papers providing some sorts of convergence guarantee Li et al. (2019b).

• The assumptions over the step-size: The three assumptions we enforced over the step-size could be easily satisfied as explained in the sequel. For instance, one can pick $\eta^t = \kappa_1 \frac{1}{t}$ for some constant κ_1 such that $\eta^t \in (0, \frac{1}{2L}]$ is satisfied. Then even if σ_t has a extremely loose upper-bound, let's say $\sigma_t < \frac{\kappa_2}{t^{\epsilon}}$ for a small $\epsilon \in \mathbb{R}_+$ and a constant number κ_2 , then all the three assumptions over the step-size in the theorem will be satisfied. Note that the convergence rate of DQN-Fed depends on how fast σ_t diminishes which depends on how heterogeneous the users are.

862 B.2 CASE 2: e > 1 & local GD

The notations used in this subsection are elaborated in Table 6.

Table 6: Notations used in the Theorem B.2 for e > 1 and local GD.

-		
	Notation	Description
-	$\boldsymbol{\theta}_{(k,e)^t}$	Updated weight for client k after e local epochs at the t -th round of FL.
-	$\mathfrak{g}_{k,e}$	$ \mathfrak{g}_{k,e} = \boldsymbol{\theta}^t - \boldsymbol{\theta}_{(k,e)^t};$ that is, the update vector of client k after e local epochs.
-	${\mathfrak g}_e$	Matrix of update vectors $[\mathfrak{g}_{1,e},\ldots,\mathfrak{g}_{K,e}]$.
-	$ ilde{\mathfrak{g}}_{k,e}$	Update vector of client k after orthogonalization process.
	$ ilde{\mathfrak{g}}_e$	Matrix of orthogonalized update vectors $[\tilde{\mathfrak{g}}_{1,e},\ldots,\tilde{\mathfrak{g}}_{K,e}]$.
-	$\lambda_{k,e}^*$	Optimum weights obtained from Equation equation 12 using $\tilde{\mathbf{g}}_e$.
-	\mathfrak{d}_e	Optimum direction obtained using $\tilde{\mathfrak{g}}_e$; that is, $\mathfrak{d}_e = \sum_{k=1}^K \lambda_{k,e}^* \tilde{\mathfrak{g}}_{k,e}$.

Theorem B.2. Assume that $f = \{f_k\}_{k \in [K]}$ are l-Lipschitz continuous and L-Lipschitz smooth. Denote by η^t and η the global and local learning rate, respectively. Also, define $\zeta^t = \|\lambda^* - \lambda_e^*\|$, where λ_e^* is the optimum weights obtained from pseudo-gradients after *e* local epochs. Then,

$$\lim_{T \to \infty} \min_{t=0,\dots,T} \|\mathbf{d}^t\| \to 0, \tag{39}$$

if the following conditions are satisfied: (i) $\eta^t \in (0, \frac{1}{2L}]$, (ii) $\lim_{T\to\infty} \sum_{t=0}^T \eta^t \to \infty$ and (iii) $\lim_{t\to\infty} \eta^t \to 0$, (iv) $\lim_{t\to\infty} \eta \to 0$, and (v) $\lim_{t\to\infty} \zeta^t \to 0$.

Proof. As discussed in the proof of Theorem B.1, we can write

$$\exists \{\lambda'_k\}_{k \in [K]} \text{ s.t. } \mathfrak{d} = \sum_{k=1}^K \lambda_k^* \tilde{\mathfrak{g}}_k = \sum_{k=1}^K \lambda'_k \mathfrak{g}_k = \mathfrak{g} \lambda',$$
(40)

To prove Theorem B.2, we first introduce a lemma whose proof is provided in Appendix C.

Lemma B.3. Using the notations used in Theorem B.2, and assumming that $f = \{f_k\}_{k \in [K]}$ are L-Lipschitz smooth, we have $\|g_{k,e} - g_k\| \le \eta el$.

 $\exists \{\lambda'_{k,e}\}_{k \in [K]} \text{ s.t. } \mathfrak{d}_e = \sum_{l=1}^{K} \lambda^*_{k,e} \tilde{\mathfrak{g}}_{k,e} = \sum_{l=1}^{K} \lambda'_{k,e} \mathfrak{g}_{k,e} = \mathfrak{g}_e \lambda'_e.$

902 Using Lemma B.3, we have

$$\|\boldsymbol{\mathfrak{d}} - \boldsymbol{\mathfrak{d}}_{e}\| = \|\tilde{\boldsymbol{\mathfrak{g}}}\boldsymbol{\lambda}^{*} - \tilde{\boldsymbol{\mathfrak{g}}}_{e}\boldsymbol{\lambda}_{e}^{*}\| \le \|\tilde{\boldsymbol{\mathfrak{g}}}\boldsymbol{\lambda}^{*} - \tilde{\boldsymbol{\mathfrak{g}}}\boldsymbol{\lambda}_{e}^{*}\| + \|\tilde{\boldsymbol{\mathfrak{g}}}\boldsymbol{\lambda}_{e}^{*} - \tilde{\boldsymbol{\mathfrak{g}}}_{e}\boldsymbol{\lambda}_{e}^{*}\|$$
(42)

$$\leq \|\tilde{\mathfrak{g}}\| \| \boldsymbol{\lambda}^* - \boldsymbol{\lambda}_e^* \| + \| \mathfrak{g} \boldsymbol{\lambda}_e' - \mathfrak{g}_e \boldsymbol{\lambda}_e' \|$$
(43)

$$\leq \|\tilde{\mathfrak{g}}\| \| \boldsymbol{\lambda}^* - \boldsymbol{\lambda}_e^* \| + \eta el \tag{44}$$

(41)

$$\zeta^t l \sqrt{K} + \eta e l, \tag{45}$$

where Equation (42) follows triangular inequality, Equation (43) is obtained from Equations (40)and (41), and Equation (44) uses Lemma B.3.

 \leq

912 As seen, if $\lim_{t\to\infty} \eta \to 0$, and $\lim_{t\to\infty} \zeta^t \to 0$, then $\|\mathbf{d} - \mathbf{d}_e\| \to 0$. Now, by writing the quadratic 913 upper bound we obtain:

$$\boldsymbol{f}(\boldsymbol{\theta}^{t+1}) \leq \boldsymbol{f}(\boldsymbol{\theta}^{t}) - \eta^{t} \boldsymbol{\mathfrak{g}}^{T} \boldsymbol{\mathfrak{g}}_{e}^{T} \boldsymbol{\lambda}_{e}^{\prime} + \frac{L(\eta^{t})^{2}}{2} \|\boldsymbol{\mathfrak{g}}_{e}^{T} \boldsymbol{\lambda}_{e}^{\prime}\|^{2}$$

$$\tag{46}$$

$$\leq \boldsymbol{f}(\boldsymbol{\theta}^{t}) - \eta^{t} \boldsymbol{\mathfrak{g}}^{T} \boldsymbol{\mathfrak{g}}^{T} \boldsymbol{\lambda}' + L(\eta^{t})^{2} \|\boldsymbol{\mathfrak{g}}^{T} \boldsymbol{\lambda}'\|^{2} + \eta^{t} \boldsymbol{\mathfrak{g}}^{T} (\boldsymbol{\mathfrak{d}} - \boldsymbol{\mathfrak{d}}_{\boldsymbol{e}}) + L(\eta^{t})^{2} \|\boldsymbol{\mathfrak{d}} - \boldsymbol{\mathfrak{d}}_{\boldsymbol{e}}\|^{2}$$
(47)

$$\leq \boldsymbol{f}(\boldsymbol{\theta}^{t}) - \eta^{t}(1 - L\eta^{t}) \|\boldsymbol{\mathfrak{g}}^{T}\boldsymbol{\lambda}'\|^{2} + l\eta^{t}\|\boldsymbol{\mathfrak{d}} - \boldsymbol{\mathfrak{d}}_{\boldsymbol{e}}\| + L(\eta^{t})^{2}\|\boldsymbol{\mathfrak{d}} - \boldsymbol{\mathfrak{d}}_{\boldsymbol{e}}\|^{2}.$$
(48)

Noting that $\eta^t \in (0, \frac{1}{2L}]$, and utilizing telescoping yields

$$\min_{t=0,\dots,T} \|\boldsymbol{\mathfrak{d}}_t\| \leq \frac{\boldsymbol{f}(\boldsymbol{\theta}_0) - \boldsymbol{f}(\boldsymbol{\theta}^{t+1}) + \sum_{t=0}^T \eta^t (l\|\boldsymbol{\mathfrak{d}} - \boldsymbol{\mathfrak{d}}_{\boldsymbol{e}}\| + L\eta^t \|\boldsymbol{\mathfrak{d}} - \boldsymbol{\mathfrak{d}}_{\boldsymbol{e}}\|^2)}{\frac{1}{2} \sum_{t=0}^T \eta^t}.$$
(49)

Using $\|\boldsymbol{\vartheta} - \boldsymbol{\vartheta}_{\boldsymbol{e}}\| \to 0$, the Theorem B.2 is concluded.

B.3 CASE 3:
$$e = 1$$
 & local GD

 Denote by ϑ the Pareto-stationary solution set of minimization problem $\arg \min_{\theta} f(\theta)$. Then, define $\boldsymbol{\theta}^* = \arg\min_{\boldsymbol{\theta}\in\vartheta} \|\boldsymbol{\theta}^t - \boldsymbol{\theta}\|_2^2$

Theorem B.4. Assume that $f = {f_k}_{k \in [K]}$ are l-Lipschitz continuous and σ -convex, and that the step-size η^t satisfies the following two conditions: (i) $\lim_{t\to\infty} \sum_{j=0}^t \eta_j \to \infty$ and (ii) $\lim_{t\to\infty}\sum_{j=0}^t\eta_j^2<\infty$. Then almost surely ${m heta}^t o {m heta}^*$; that is,

$$\mathbb{P}\left(\lim_{t\to\infty}\left(\boldsymbol{\theta}^t - \boldsymbol{\theta}^*\right) = 0\right) = 1,\tag{50}$$

where $\mathbb{P}(E)$ denotes the probability of event E.

Proof. The proof is inspired from Mercier et al. (2018). Without loss of generality, we assume that all users participate in all rounds.

Based on the definition of θ^* we can say

$$\|\boldsymbol{\theta}^{t+1} - \boldsymbol{\theta}_{t+1}^*\|_2^2 \le \|\boldsymbol{\theta}^{t+1} - \boldsymbol{\theta}_t^*\|_2^2 = \|\boldsymbol{\theta}^t - \eta^t \mathfrak{d}_t - \boldsymbol{\theta}_t^*\|_2^2$$
(51)

$$= \|\boldsymbol{\theta}^t - \boldsymbol{\theta}_t^*\|_2^2 - 2\eta^t (\boldsymbol{\theta}^t - \boldsymbol{\theta}_t^*) \cdot \boldsymbol{\mathfrak{d}}_t + (\eta^t)^2 \|\boldsymbol{\mathfrak{d}}_t\|_2^2.$$
(52)

= To bound the third term in Equation (52), we note that from Equation (27), we have:

$$(\eta^t)^2 \|\mathbf{d}_t\|_2^2 = \frac{(\eta^t)^2}{\sum_{k=1}^K \frac{1}{\|\bar{\mathbf{g}}_k\|_2^2}} \le \frac{(\eta^t)^2 l^2}{K}.$$
(53)

To bound the second term, first note that since orthogonal vectors $\{\tilde{\mathfrak{g}}_k\}_{k\in[K]}$ span the same Kdimensional space as that spanned by gradient vectors $\{\mathfrak{g}_k\}_{k\in[K]}$, then

$$\exists \{\lambda'_k\}_{k \in [K]} \text{ s.t. } \mathfrak{d} = \sum_{k=1}^K \lambda^*_k \tilde{\mathfrak{g}}_k = \sum_{k=1}^K \lambda'_k \mathfrak{g}_k.$$
(54)

Using Equation (54) and the σ -convexity of $\{f_k\}_{k \in [K]}$ we obtain

$$(\boldsymbol{\theta}^{t} - \boldsymbol{\theta}_{t}^{*}) \cdot \boldsymbol{\mathfrak{d}}_{t} = (\boldsymbol{\theta}^{t} - \boldsymbol{\theta}_{t}^{*}) \cdot \sum_{k=1}^{K} \lambda_{k}^{*} \tilde{\boldsymbol{\mathfrak{g}}}_{k}$$
(55)

$$= (\boldsymbol{\theta}^{t} - \boldsymbol{\theta}_{t}^{*}) \cdot \sum_{k=1}^{K} \lambda_{k}^{\prime} \boldsymbol{\mathfrak{g}}_{k}$$
(56)

$$\geq \sum_{k=1}^{K} \lambda_k' \left(f_k(\boldsymbol{\theta}^t) - f_k(\boldsymbol{\theta}_t^*) \right) + \sigma \frac{\|\boldsymbol{\theta}^t - \boldsymbol{\theta}_t^*\|_2^2}{2}$$
(57)

$$\geq \frac{\lambda_{\alpha}'M}{2} \|\boldsymbol{\theta}^t - \boldsymbol{\theta}_t^*\|_2^2 + \sigma \frac{\|\boldsymbol{\theta}^t - \boldsymbol{\theta}_t^*\|_2^2}{2}$$
(58)

 $=\frac{\lambda'_{\alpha}M+\sigma}{2}\|\boldsymbol{\theta}^t-\boldsymbol{\theta}^*_t\|_2^2.$ (59)

Now, we return back to Equation (52) and find the conditional expectation w.r.t. θ^t as follows

 $\mathbf{E}[\|\boldsymbol{\theta}^{t+1} - \boldsymbol{\theta}_{t+1}^*\|_2^2 \mid \boldsymbol{\theta}^t] \le (1 - \eta^t \mathbf{E}[\lambda_{\alpha}'M + \sigma|\boldsymbol{\theta}^t])\|\boldsymbol{\theta}^t - \boldsymbol{\theta}_t^*\|_2^2 + \frac{(\eta^t)^2 l^2}{\kappa}.$ (60)

Assume that $\mathbf{E}[\lambda'_{\alpha}M + \sigma | \boldsymbol{\theta}^t] \geq c$, taking another expectation we obtain:

$$\mathbf{E}[\|\boldsymbol{\theta}^{t+1} - \boldsymbol{\theta}_{t+1}^*\|_2^2] \le (1 - \eta^t c) \mathbf{E}[\|\boldsymbol{\theta}^t - \boldsymbol{\theta}_t^*\|_2^2] + \frac{(\eta^t)^2 l^2}{K},\tag{61}$$

which is a recursive expression. By solving Equation (61) we obtain

$$\mathbf{E}[\|\boldsymbol{\theta}^{t+1} - \boldsymbol{\theta}^{*}_{t+1}\|_{2}^{2}] \leq \underbrace{\prod_{j=0}^{t} (1 - \eta_{j}c) \mathbf{E}[\|\boldsymbol{\theta}_{0} - \boldsymbol{\theta}^{*}_{0}\|_{2}^{2}]}_{\text{First term}} + \underbrace{\sum_{m=1}^{t} \frac{\prod_{j=1}^{t} (1 - \eta_{j}c) \eta_{m}^{2} l^{2}}{K \prod_{j=1}^{m} (1 - \eta_{j}c)}}_{\text{Second term}}.$$
 (62)

It is observed that if the limit of both First term and Second term in Equation (62) go to zero, then $\mathbf{E}[\|\boldsymbol{\theta}^{t+1} - \boldsymbol{\theta}^*_{t+1}\|_2^2] \to 0.$ For the First term, from the arithmetic-geometric mean inequality we have

$$\lim_{t \to \infty} \prod_{j=0}^{t} (1 - \eta_j c) \le \lim_{t \to \infty} \left(\frac{\sum_{j=0}^{t} (1 - \eta_j c)}{t} \right)^t = \lim_{t \to \infty} \left(1 - c \frac{\sum_{j=0}^{t} \eta_j}{t} \right)^t \tag{63}$$

$$=\lim_{t\to\infty}e^{-c\sum_{j=0}^t\eta_j}.$$
(64)

From Equation (64) it is seen that if $\lim_{t\to\infty} \sum_{j=0}^t \eta_j \to \infty$, then the First term is also converges to zero as $t \to \infty$.

On the other hand, consider the Second term in Equation (62). Obviously, if $\lim_{t\to\infty} \sum_{j=0}^t \eta_j^2 < \infty$, then the Second term converges to zero as $t \to \infty$.

Hence, if (i) $\lim_{t\to\infty} \sum_{j=0}^t \eta_j \to \infty$ and (ii) $\lim_{t\to\infty} \sum_{j=0}^t \eta_j^2 < \infty$, then $\mathbf{E}[\|\boldsymbol{\theta}^{t+1} - \boldsymbol{\theta}_{t+1}^*\|_2^2] \to 0$. Consequently, based on standard supermartingale Mercier et al. (2018), we have

$$\mathbb{P}\left(\lim_{t\to\infty}\left(\boldsymbol{\theta}^t - \boldsymbol{\theta}^*\right) = 0\right) = 1.$$
(65)

PROOF OF LEMMA B.3 С

Proof.

$$\mathbf{g}_{k,e} = \mathbf{\theta}^{t} - \mathbf{\theta}_{(k,e)^{t}} = (\mathbf{\theta}^{t} - \mathbf{\theta}_{(k,1)^{t}}) + (\mathbf{\theta}_{(k,1)^{t}} - \mathbf{\theta}_{(k,2)^{t}}) + \dots + (\mathbf{\theta}_{(k,e-1)^{t}} - \mathbf{\theta}_{(k,e)^{t}})$$
(66)

$$=\mathfrak{g}_k(\boldsymbol{\theta}^t) + \eta\mathfrak{g}_{k,1} + \dots + \eta\mathfrak{g}_{k,e-1}.$$
(67)

Hence,

$$\|\mathfrak{g}_{k,e} - \mathfrak{g}_k\| = \|\eta \sum_{j=1}^e \mathfrak{g}_{k,j}\| \le \eta \sum_{j=1}^e \|\mathfrak{g}_{k,j}\| \le \eta el.$$
(68)

D **CONVERGENCE RATE, FULL-VERSION OF THEOREM 4.6**

In this subsection, we provide convergence guarantee for DQN-Fed. First, consider the following assumptions which we use in our convergence theorems.

Definition D.1. We say that \mathbf{B}_t^{-1} is a δ -approximate of true inverse Hessian \mathbf{H}_t^{-1} , if the following holds

$$\|\mathbf{B}_{t}^{-1} - \mathbf{H}_{t}^{-1}\| \le \delta \|\mathbf{H}_{t}^{-1}\|.$$
(69)

• Assumption 1. The global loss function is twice continuously differentiable, L-Lipschitz gradient (*L*-smooth) and λ -strongly convex. As such, we have

$$\lambda \mathbf{I} \preceq \nabla^2 f(\boldsymbol{\theta}) \preceq L \mathbf{I},\tag{70}$$

where I represents the identity matrix, and the notation $A \leq B$, where A and B are matrices of the same size, signifies that $\mathbf{A} - \mathbf{B}$ is positive semidefinite. Note that the strong convexity of the global loss function implies that there exists a unique optimal model parameter which we denote by θ^{\star} hereafter in our proof.

• Assumption 2. The matrix \mathbf{B}_t^{-1} is a δ -approximate of true inverse Hessian \mathbf{H}_t^{-1} .

Theorem D.2. Let Assumptions 1 & 2 hold. Then,

$$\|\boldsymbol{\theta}_{t} - \boldsymbol{\theta}^{\star}\| \leq \begin{cases} \left(\frac{L\delta}{\lambda}\right)^{t} \|\boldsymbol{\theta}_{0} - \boldsymbol{\theta}^{\star}\| + A_{0}^{\prime}, & t \leq t_{0} \\ \left(\frac{L\delta}{\lambda}\right)^{t} \|\boldsymbol{\theta}_{0} - \boldsymbol{\theta}^{\star}\| + A_{1}^{\prime}, & t > t_{0} \end{cases}$$
(71)

where A'_0 and A'_1 are defined as follows

$$A_0' = \frac{\left(\frac{L\delta}{\lambda}\right)^t - 1}{\frac{L\delta}{\lambda} - 1} \left[\frac{\lambda}{L} \left(t_0 - t + \frac{2\gamma}{1 - \gamma}\right)\right],\tag{72a}$$

$$A_1' = \frac{\left(\frac{L\delta}{\lambda}\right)^t - 1}{\frac{L\delta}{\lambda} - 1} \left[\frac{2\lambda\gamma^{2^{t-t_0}}}{L(1 - \gamma^{2^{t-t_0}})} + \right],\tag{72b}$$

with
$$t_0 = \max\left\{0, \left\lceil\frac{2L}{\lambda^2 \|\boldsymbol{\mathfrak{d}}_0\|}\right\rceil - 2\right\},$$
 (72c)

and
$$\gamma = \frac{L}{2\lambda^2} \|\boldsymbol{\vartheta}_0\| - \frac{t_0}{4}.$$
 (72d)

Proof. The proof is differed to Appendix E.

As per Theorem equation D.2, DQN-Fed method has a linear-quadratic convergence rate. In fact, the quadratic term in equation 71 is exactly the same as that of Polyak & Tremba (2020); yet, the linear term is the result of approximating the local Hessian matrices using BFGS method.

In the ensuing corollary, our objective is to determine the required number of communication rounds T_{ϵ} such that $\|\boldsymbol{\theta}_{T_{\epsilon}} - \boldsymbol{\theta}^{\star}\| \leq \epsilon$.

Corollary D.3. If
$$\|\boldsymbol{\theta}_0 - \boldsymbol{\theta}^{\star}\| < \frac{A'_1}{\left(\frac{L\delta}{\lambda}\right)^t}$$
, then DQN-Fed has a quadratic convergence rate:

$$\|\boldsymbol{\theta}_t - \boldsymbol{\theta}^\star\| \le 2A_1'. \tag{73}$$

Also, if $\frac{L\delta}{\lambda} < 1$, we have

> $T_{\epsilon} = \mathcal{O}\left(\log\log\frac{1}{\epsilon}\right),$ (74)

which is also called super-linear convergence rate.

Corollary D.4. On the other hand, If $\left\| \boldsymbol{\theta}_0 - \boldsymbol{\theta}^{\star} \right\| \geq \frac{A_1'}{\left(\frac{L\delta}{\lambda}\right)^t}$ and $\frac{L\delta}{\lambda} < 1$, then DQN-Fed method has a linear convergence rate:

$$\|\boldsymbol{\theta}_t - \boldsymbol{\theta}^{\star}\| \le 2 \left(\frac{L\delta}{\lambda}\right)^t \|\boldsymbol{\theta}_0 - \boldsymbol{\theta}^{\star}\|.$$
(75)

and,

$$T_{\epsilon} = \mathcal{O}\left(\frac{1}{\log(\frac{\lambda}{L\delta})}\log\frac{1}{\epsilon}\right).$$
(76)

The proof for Corollary D.3 and D.4 can be found in Appendix F and G, respectively.

Remark D.5. It is worth noting that for distributed GD-like methods, the number of communication rounds needed to achieve a desired precision ϵ , follows a linear convergence rate. Specifically, we have $T_{\epsilon} = \mathcal{O}\left(\frac{L}{\lambda}\log\frac{1}{\epsilon}\right)$. This underscores the superiority of DQN-Fed in terms of convergence rate.

¹⁰⁸⁰ E PROOF OF THEOREM THEOREM D.2

Throughout the proofs in this section, we frequently use the triangular inequality for two vectors **v** and **u**: $\|\mathbf{v} \pm \mathbf{u}\| \le \|\mathbf{v}\| + \|\mathbf{u}\|$.

Our goal is to prove the theorem by deriving a recursive relation for the distance between the optimal global model θ^* and the global model at the *t*-th round θ_t , specifically $\|\theta_t - \theta^*\|$. First, noting that $\theta_{t+1} = \theta_t + \eta_t \mathbf{B}_t^{-1} \tilde{\mathbf{g}}_t$, we have

 $\|\boldsymbol{\theta}_{t+1} - \boldsymbol{\theta}^{\star}\| = \left\| \boldsymbol{\theta}_{t} - \eta_{t} \mathbf{B}_{t}^{-1} \tilde{\mathbf{g}}_{t} - \boldsymbol{\theta}^{\star} \right\| \leq \underbrace{\left\| \boldsymbol{\theta}_{t} - \eta_{t} \mathbf{H}_{t}^{-1} \mathbf{g}_{t} - \boldsymbol{\theta}^{\star} \right\|}_{M_{1}} + \eta_{t} \underbrace{\left\| \mathbf{H}_{t}^{-1} \mathbf{g}_{t} - \mathbf{B}_{t}^{-1} \tilde{\mathbf{g}}_{t} \right\|}_{M_{2}}.$ (77)

1095 To bound M_1 , we use the results in Polyak & Tremba (2020). In particular, define $t_0 = \max\left\{0, \left\lceil\frac{2L}{\lambda^2 \|\boldsymbol{\mathfrak{d}}^0\|}\right\rceil - 2\right\}, \gamma = \frac{L}{2\lambda^2} \|\boldsymbol{\mathfrak{d}}^0\| - \frac{t_0}{4}$; then, 1097

1098
1099
1100
$$M_{1} \leq \begin{cases} \frac{\lambda}{L}(t_{0} - t + \frac{2\gamma}{1-\gamma}), & t \leq t_{0} \\ \frac{2\lambda\gamma^{2^{t-t_{0}}}}{L(1-\gamma^{2^{t-t_{0}}})}, & t > t_{0} \end{cases}$$
(78)

1101 Next, we bound M_2 in the sequel.

1088

1089 1090 1091

1093 1094

1103 1104 1105

1109

1110

1112 1113

1122 1123 1124

1126 1127 1128

$$M_2 \le \left\| \mathbf{H}_t^{-1} \mathbf{g}_t - \mathbf{B}_t^{-1} \mathbf{g}_t \right\| + \left\| \mathbf{B}_t^{-1} \mathbf{g}_t - \mathbf{B}_t^{-1} \tilde{\mathbf{g}}_t \right\|$$
(79)

$$\leq \left\|\mathbf{H}_{t}^{-1} - \mathbf{B}_{t}^{-1}\right\| \left\|\mathbf{g}_{t}\right\| + \left\|\mathbf{B}_{t}^{-1}\right\| \left\|\mathbf{g}_{t} - \tilde{\mathbf{g}}_{t}\right\|,\tag{80}$$

where in equation 79 we used triangular inequality. Note that using the assumption 2, we have $\|\mathbf{H}_t^{-1} - \mathbf{B}_t^{-1}\| \le \delta \|\mathbf{H}_t^{-1}\|$, and by λ -strong convexity of the loss function, we have $\|\mathbf{H}_t^{-1}\| \le \frac{1}{\lambda}$. Hence,

$$\left|\mathbf{H}_{t}^{-1} - \mathbf{B}_{t}^{-1}\right| \leq \frac{\delta}{\lambda}.$$
(81)

¹¹¹¹ In addition, the *L*-smoothness of the global loss function yields

$$\left\|\mathbf{g}_{t}\right\| \leq L\left\|\boldsymbol{\theta}_{t} - \boldsymbol{\theta}^{\star}\right\|.$$
(82)

Hence, from equation 81 and equation 82, the first term in equation 80 could be bounded. Now, tobound the second term in equation 80, note that

$$\mathbf{B}_{t}^{-1} \| \leq \left\| \mathbf{B}_{t}^{-1} - \mathbf{H}_{t}^{-1} \right\| + \left\| \mathbf{H}_{t}^{-1} \right\|$$
(83)

$$\leq \frac{\delta}{\lambda} + \frac{1}{\lambda} = \frac{\delta + 1}{\lambda}.$$
(84)

Using equation 81, equation 82 and equation 83 in the inequality equation 79 we obtain

$$M_2 \le \frac{L\delta}{\lambda} \|\boldsymbol{\theta}_t - \boldsymbol{\theta}^\star\| + \frac{\delta+1}{\lambda} \|\mathbf{g}_t - \tilde{\mathbf{g}}_t\|.$$
(85)

1125 Next, we have (note that $\eta_t \leq 1$)

$$\|\boldsymbol{\theta}_{t+1} - \boldsymbol{\theta}^{\star}\| \leq \begin{cases} \frac{L\delta}{\lambda} \|\boldsymbol{\theta}_t - \boldsymbol{\theta}^{\star}\| + A_0, & t \leq t_0\\ \frac{L\delta}{\lambda} \|\boldsymbol{\theta}_t - \boldsymbol{\theta}^{\star}\| + A_1, & t > t_0 \end{cases}$$
(86a)

where
$$A_0 = \frac{\lambda}{L} (t_0 - t + \frac{2\gamma}{1 - \gamma}),$$
 (86b)

and
$$A_1 = \frac{2\lambda\gamma^{2^{t-10}}}{L(1-\gamma^{2^{t-t_0}})}.$$
 (86c)

Applying equation 86 recursively yields

$$\|\boldsymbol{\theta}_{t} - \boldsymbol{\theta}^{\star}\| \leq \begin{cases} \left(\frac{L\delta}{\lambda}\right)^{t} \|\boldsymbol{\theta}_{0} - \boldsymbol{\theta}^{\star}\| + A_{0}^{\prime}, & t \leq t_{0} \\ \left(\frac{L\delta}{\lambda}\right)^{t} \|\boldsymbol{\theta}_{0} - \boldsymbol{\theta}^{\star}\| + A_{1}^{\prime}, & t > t_{0} \end{cases}$$
(87a)

(87b)

(88)

where
$$A_0' = rac{\left(rac{L\delta}{\lambda}
ight)^t - 1}{rac{L\delta}{\lambda} - 1}A_0,$$

and
$$A'_{1} = \frac{\left(\frac{L\delta}{\lambda}\right)^{t} - 1}{\frac{L\delta}{\lambda} - 1} A_{1}.$$
 (87c)

¹¹⁴⁴ ₁₁₄₅ F Proof of Corollary D.3

1147 As per Theorem equation D.2, if $\|\boldsymbol{\theta}_0 - \boldsymbol{\theta}^\star\| < \frac{A_1'}{\left(\frac{L\delta}{\lambda}\right)^t}$, then 1149 $\|\boldsymbol{\theta}_t - \boldsymbol{\theta}^\star\| \le \left(\frac{L\delta}{\lambda}\right)^t \frac{A_1'}{\left(\frac{L\delta}{\lambda}\right)^t} + A_1' = 2A_1'.$

1151 Hence, to find T_{ϵ} , we shall have

$$2A_1' \le \epsilon \tag{89}$$

$$\Leftrightarrow 2 \frac{\left(\frac{L\delta}{\lambda}\right)^{t} - 1}{\frac{L\delta}{\lambda} - 1} \left[\underbrace{\frac{2\lambda\gamma^{2^{t-t_0}}}{L(1 - \gamma^{2^{t-t_0}})}}_{\text{diminishing term}} \right] \le \epsilon.$$
(90)

1157 (158)
1157 Since
$$\left(\frac{L\delta}{\lambda}\right) < 1$$
, as t becomes larger, $\left(\frac{L\delta}{\lambda}\right)^t \approx 0$, and therefore $2\frac{\left(\frac{L\delta}{\lambda}\right)^t - 1}{\frac{L\delta}{\lambda} - 1} \approx \frac{2}{1 - \frac{L\delta}{\lambda}}$. In addition,
1159 since $\gamma \in [0, \frac{1}{2}]$, for the large values of t, $\frac{2\lambda\gamma^{2^{t-t_0}}}{L(1-\gamma^{2^{t-t_0}})} \approx \frac{2\lambda\gamma^{2^{t-t_0}}}{L}$. Thus, by inverting the inequality
1161 equation 90, and then taking log from both sides we have

$$\log(\frac{1}{\epsilon}) \le -\log(\frac{4\lambda}{L - \frac{L^2\delta}{\lambda}}) - 2^{t-t_0}\log(\gamma).$$
(91)

1165 Note that $\log(\gamma) < 0$, and therefore the second term on the RHS of equation 91 is positive. Also, 1166 since $-\log(\frac{4\lambda}{L-\frac{L^2\delta}{\lambda}}) \ll -2^{t-t_0}\log(\gamma)$, then

$$T_{\epsilon} \leq \mathcal{O}\left(\log\log\frac{1}{\epsilon}\right).$$
 (92)

G PROOF OF COROLLARY D.4

Based on Theorem equation D.2, if $\|\boldsymbol{\theta}_0 - \boldsymbol{\theta}^\star\| \ge \frac{A_1'}{\left(\frac{L\delta}{\lambda}\right)^t}$, then

 \leq

$$\|\boldsymbol{\theta}_{t} - \boldsymbol{\theta}^{\star}\| \leq \left(\frac{L\delta}{\lambda}\right)^{t} \|\boldsymbol{\theta}_{0} - \boldsymbol{\theta}^{\star}\| + \left(\frac{L\delta}{\lambda}\right)^{t} \|\boldsymbol{\theta}_{0} - \boldsymbol{\theta}^{\star}\|$$
(93)

$$2\left(\frac{L\delta}{\lambda}\right)^{t} \left\|\boldsymbol{\theta}_{0}-\boldsymbol{\theta}^{\star}\right\|.$$
(94)

1179 Thus, to find T_{ϵ} , we shall have

$$2\left(\frac{L\delta}{\lambda}\right)^{T_{\epsilon}} \left\|\boldsymbol{\theta}_{0} - \boldsymbol{\theta}^{\star}\right\| \leq \epsilon$$
(95)

$$\Leftrightarrow \qquad T_{\epsilon} \log(\frac{\lambda}{L\delta}) \ge \log\left(\frac{2\|\boldsymbol{\theta}_0 - \boldsymbol{\theta}^{\star}\|}{\epsilon}\right). \tag{96}$$

$$T_{\epsilon} = \mathcal{O}\left(\frac{1}{\log(\frac{\lambda}{L\delta})}\log\frac{1}{\epsilon}\right).$$
(97)

1188 H EXPERIMENTAL ANALYSIS

1190 H.1 COMPARISON OF CONVERGENCE RATE

1192 In this subsection, we empirically compare the convergence speed of DQN-Fed against several fair 1193 first-order methods. To do so, we use the four datasets from setup 1 described in Section 5 and plot 1194 the validation accuracy of different methods as a function of communication rounds. The results are 1195 shown in Figure 1. As observed, DQN-Fed demonstrates a faster convergence rate compared to all 1196 benchmark methods.





1221 H.1.1 PERCENTAGE OF IMPROVED CLIENTS

We measure the training loss before and after each communication round for all participating clients and report the percentage of clients whose loss function decreased or remained unchanged, defined $\frac{\sum_{k \in S_t} \mathbb{I}\{f_k(\theta^{t+1}) \leq f_k(\theta^t)\}}{|S|}, \text{ where } S_t \text{ is the participating clients in round } t, \text{ and } \mathbb{I}(\cdot) \text{ is the } t$ as $\rho_t =$ indicator function. Then, we plot ρ_t versus communication rounds for different fair FL methods. The curves for CIFAR-10 and CIFAR-100 datasets are reported in Figure 2a and Figure 2b, respectively. As seen, both DQN-Fed and FedMGDA+ consistently outperform other benchmark methods in that fewer clients' performances get worse after participation. We further note that after enough number of communication rounds, curves for both DQN-Fed and FedMGDA+ converge to 100%. In ??, we also present experimental analysis to validate the claims made in the paper.

I ADDITIONAL DATASETS

In this section, we assess the performance of DQN-Fed against several benchmarks using additional
 datasets, namely Fashion MNIST, CINIC-10, and TinyImageNet. The corresponding results for each
 dataset are detailed in Appendices I.1 to I.3.

1239 I.1 FASHION MNIST

Fashion MNIST (Xiao et al., 2017) is an extension of MNIST dataset (LeCun et al., 1998) with images resized to 32×32 pixels.



We use a fully-connected neural network with 2 hidden layers, and use the same setting as that used in Li et al. (2019a) for our experiments. We set e = 1 and use full batchsize, and use $\eta = 0.1$. Then, we conduct 300 rounds of communications. For the benchmarks, we use the same as those we used for CIFAR-10 experiments. The results are reported in Table 7.

By observing the three different classes reported in Table 7, we observe that the fairness level attained in DQN-Fed is not limited to a dominate class.

Table 7: Test accuracy on Fashion MNIST. The reported results are averaged over 5 different seeds.

Algorithm	$ \bar{a}$	σ_a	SHIRT	PULLOVER	T-SHIRT
FEDAVG	80.42 78.53 79.29 80.22 81.27	3.39	64.26	87.00	89.90
Q-FFL		2.27	71.29	81.46	82.86
FEDMGDA+		2.53	<u>72.46</u>	79.74	85.66
FEDHEAL		3.41	63.71	86.87	89.94
DQN-FED		<u>2.31</u>	72.57	88.21	90.99

1276 1277 I.2 CINIC-10

1257

1265

1287

1289 1290 1291

1293 1294 1295

1278 CINIC-10 (Darlow et al., 2018) has 4.5 times as many images as those in CIFAR-10 dataset (270,000 sample images in total). In fact, it is obtained from ImageNet and CIFAR-10 datasets. As a result, this dataset fits FL scenarios since the constituent elements of CINIC-10 are not drawn from the same distribution. Furthermore, we add more non-iidness to the dataset by distributing the data among the clients using Dirichlet allocation with $\beta = 0.5$.

For the model, we use ResNet-18 with group normalization, and set $\eta = 0.01$. There are 200 communication rounds in which all the clients participate with e = 1. Also, K = 50. Results are reported in Table 8.

Table 8: Test accuracy on CINIC-10. The reported results are averaged over 5 different seeds.

Algorithm	\bar{a}	σ_a	WORST 10%	Best 10%
Q-FFL	86.57	14.91	57.70	100.00
DITTO	86.31	15.14	56.91	100.00
FedLF	86.49	15.12	57.62	100.00
TERM	86.40	15.10	57.30	100.00
DQN-FED	87.34	14.85	57.88	99.99

1296 I.3 TINYIMAGENET

1298 Tiny-ImageNet (Le & Yang, 2015) is a subset of ImageNet with 100k samples of 200 classes. We 1299 distribute the dataset among K = 20 clients using Dirichlet allocation with $\beta = 0.05$

We use ResNet-18 with group normalization, and set $\eta = 0.02$. There are 400 communication rounds in which all the clients participate with e = 1. The results are reported in Table 9.

1303Table 9: Test accuracy on TinyImageNet. The reported results are averaged over 5 different seeds.

Algorithm	ā	σ_a	WORST 10%	Best 10%
Q-FFL	18.90	3.20	13.12	23.72
FedLF	16.55	2.38	12.40	20.25
TERM	16.41	2.77	11.52	21.02
FedMGDA+	14.00	2.71	9.88	19.21
DQN-FED	19.05	2.35	13.24	23.58

1310 1311 1312

1313 1314

1316 1317

1318

1319

1320 1321

1322

1323 1324

1325 1326

1327

1309

1304 1305 1306

J EXPERIMENTS DETAILS, TUNING HYPER-PARAMETERS

For all benchmark methods, we conducted a grid-search to identify the optimal hyper-parameters for the underlying algorithms. The parameters tested for each method are outlined below:

• **q-FFL:** $q \in \{0, 0.001, 0.01, 0.1, 1, 2, 5, 10\}.$

• **TERM:** $t \in \{0.1, 0.5, 1, 2, 5\}.$

• FedLF: $\eta^t \in \{0.01, 0.05, 0.1, 0.5, 1\}.$

• **Ditto:** $\lambda \in \{0.01, 0.05, 0.1, 0.5, 1, 2, 5\}.$

• FedMGDA+: $\epsilon \in \{0.01, 0.05, 0.1, 0.5, 1\}.$

• FedHEAL: $(\alpha, \beta) = \{(0.5, 0.5)\}, (\gamma_s, \gamma_c) = \{(0.5, 0.9)\}.$

K INTEGRATION WITH A LABEL NOISE CORRECTION METHOD

1328 1329 K.1 Understanding Label Noise in FL

Label noise in FL refers to inaccuracies or errors present in the ground truth labels associated with
the data used for training FL models. These inaccuracies manifest when the labels assigned to data
points are incorrect or noisy due to various reasons. Label noise can originate at different stages
of the FL process, including data collection, annotation, or transmission phases. Addressing label
noise is crucial as it can substantially impact the performance and reliability of FL models, making it
essential to develop robust strategies to mitigate its effects.

Addressing label noise in FL presents unique challenges due to its reliance on decentralized data
 sources, where participants may have limited control over label quality in remote environments.
 Mitigating label noise in FL requires the development of robust models and FL algorithms capable of
 adapting to inaccuracies in the labels. This adaptation is essential for maintaining model performance
 and reliability in real-world FL scenarios where label noise is prevalent.

1341

1342 K.2 ROBUSTNESS OF FAIR FL ALGORITHMS TO LABEL NOISE

The core objective of fair FL algorithms, such as DQN-Fed, is to uphold fairness among clients while
preserving average accuracy across them. However, it's important to note that these algorithms are
not inherently robust against label noise, which refers to instances where data points are mislabeled.

However, by integrating DQN-Fed with label-noise resistant methods from existing literature, we
can develop a FL approach that not only ensures fairness among clients but also exhibits robustness
against label noise. Specifically, among the label-noise resistant FL algorithms available in the
literature, we choose FedCorr (Xu et al., 2022) to be integrated with DQN-Fed. This integration

offers a promising avenue for enhancing the performance and resilience of FL models in real-world
 scenarios affected by label noise.

FedCorr introduces a dimensionality-based filter to identify noisy clients, achieved through the measurement of local intrinsic dimensionality (LID) within local model prediction subspaces. They illustrate the feasibility of distinguishing between clean and noisy datasets by monitoring the behavior of LID scores throughout the training process. For further insights into FedCorr, we defer interested readers to the original paper for a comprehensive discussion.

Following a methodology similar to FedCorr, we utilize a real-world noisy dataset known as Clothing1M. This dataset comprises 1 million clothing images across 14 classes and is characterized by noisy labels, as it is sourced from various online shopping websites, incorporating numerous mislabeled samples.

For our experiments with Clothing1M, we adopt the identical settings as utilized by FedCorr, which are available in their GitHub repository (https://github.com/Xu-Jingyi/FedCorr). Specifically, we employ local SGD with a momentum of 0.5, utilizing a batch size of 16 and conducting five local epochs. Additionally, we set the hyper-parameter $T_1 = 2$ in accordance with their algorithm.

The results are summarized in Table 10. It is evident that the average accuracy achieved by DQNFed is approximately 2.2% lower compared to that obtained with FedCorr, indicating DQN-Fed's susceptibility to label noise. However, DQN-Fed demonstrates a notable improvement in ensuring fair client accuracy, aligning with expectations.

Conversely, when DQN-Fed is combined with FedCorr, there is a noticeable enhancement in average
 accuracy while still preserving satisfactory fairness among clients. This integration showcases the
 potential of leveraging both methodologies to achieve improved performance and fairness in FL
 scenarios affected by label noise.

Table 10: Test accuracy on Clothing1M dataset. The reported results are averaged over 5 different seeds.

Algorithm	$ \bar{a}$	σ_a	W(10%)	B(10%)
FEDAVG	70.49	13.25	43.09	91.05
FEDCORR	72.55	13.27	43.12	91.15
DQN-FED	70.35	5.17	49.91	90.77
FEDCORR + DON-FED	72.36	8.07	46.77	91.15

1382 1383 1384

1377

1380 1381

1385 1386

L MORE ON FAIRNESS IN FL AND ML

1387 L.1 SOURCES OF UNFAIRNESS IN FEDERATED LEARNING 1388

Unfairness in FL can arise from various sources and is a concern that needs to be addressed in FLsystems. Here are some of the key reasons for unfairness in FL:

1. Non-Representative Data Distribution: Unfairness can occur when the distribution of data across participating devices or clients is non-representative of the overall population. Some devices may have more or less relevant data, leading to biased model updates.

1394
1395
2. Data Bias: If the data collected or used by different clients is inherently biased due to the data collection process, it can lead to unfairness. For example, if certain demographic groups are underrepresented in the training data of some clients, the federated model may not perform well for those groups.

3. Heterogeneous Data Sources: Federated learning often involves data from a diverse set of sources, including different device types, locations, or user demographics. Variability in data sources can introduce unfairness as the models may not generalize equally well across all sources.

4. Varying Data Quality: Data quality can vary among clients, leading to unfairness. Some clients may have noisy or less reliable data, while others may have high-quality data, affecting the model's performance.

5. Data Sampling: The way data is sampled and used for local updates can introduce unfairness. If some clients have imbalanced or non-representative data sampling strategies, it can lead to biased model updates.

Aggregation Bias: The learned model may exhibit a bias towards devices with larger amounts of data or, if devices are weighted equally, it may favor more commonly occurring devices.

1409 1410 1411

L.2 FAIRNESS IN CONVENTIONAL ML VS. FL

The concept of fairness is often used to address social biases or performance disparities among different individuals or groups in the machine learning (ML) literature (Barocas et al., 2017). However, in the context of FL, the notion of fairness differs slightly from traditional ML. In FL, fairness primarily pertains to the consistency of performance across various clients. In fact, the difference in the notion of fairness between traditional ML and FL arises from the distinct contexts and challenges of these two settings:

1418 1419 1420

1421

1422

1423

1424

1425

1426

1. Centralized vs. decentralized data distribution:

- In traditional ML, data is typically centralized, and fairness is often defined in terms of mitigating biases or disparities within a single, homogeneous dataset. Fairness is evaluated based on how the model treats different individuals or groups within that dataset.
- In FL, data is distributed across multiple decentralized clients or devices. Each client may have its own unique data distribution, and fairness considerations extend to addressing disparities across these clients, ensuring that the federated model provides uniform and equitable performance for all clients.
- 1427 1428 1429 1430

1431

1432

1433 1434

1435 1436

2. Client autonomy and data heterogeneity:

- In FL, clients are autonomous and may have different data sources, labeling processes, and data collection practices. Fairness in this context involves adapting to the heterogeneity and diversity among clients while still achieving equitable outcomes.
 - Traditional ML operates under a centralized, unified data schema and is not inherently designed to handle data heterogeneity across sources.

We should note that in certain cases where devices can be naturally clustered into groups with specific attributes, the definition of fairness in FL can be seen as a relaxed version of that in ML, i.e., we optimize for similar but not necessarily identical performance across devices (Li et al., 2019a).

Nevertheless, despite the differences mentioned above, to maintain consistency with the terminology used in the FL literature and the papers we have cited in the main body of this work, we will continue to use the term "fairness" to denote the uniformity of performance across different devices.

1443

1445

1444 L.3 FAIR ALGORITHMS IN FL

A seminal method in this domain is Agnostic Federated Learning (FedLF) Mohri et al. (2019). 1446 FedLF optimizes the global model for the worst-case realization of the weighted combination of 1447 user distributions. Their approach involves solving a saddle-point optimization problem, and they 1448 employ a fast stochastic optimization algorithm for this purpose. However, FedLF exhibits strong 1449 performance only for a limited number of clients. In addition, Du et al. (2021) adopted the framework 1450 of FedLF and introduced the AgnosticFair algorithm. They linearly parameterized model weights 1451 using kernel functions and demonstrated that FedLF can be considered as a specific instance of 1452 AgnosticFair. To address the challenges in FedLF, the q-fair Federated Learning (q-FFL) method was 1453 introduced by Li et al. (2019a). q-FFL aims to achieve a more uniform test accuracy across users, 1454 drawing inspiration from fair resource allocation methods employed in wireless communication networks Huaizhou et al. (2013). Following this, Li et al. (2020) introduced TERM, a tilted empirical 1455 risk minimization algorithm designed to address outliers and class imbalance in statistical estimation 1456 procedures. In numerous FL applications, TERM has shown superior performance compared to 1457 q-FFL. Adopting a similar concept, Huang et al. (2020b) introduced a method that adjusts device

weights based on training accuracy and frequency to promote fairness. Additionally, FCFC Cui et al. (2021) minimizes the loss of the worst-performing client, effectively presenting a variant of FedLF. Subsequently, Li et al. (2021) introduced Ditto, a multitask personalized FL algorithm. By optimizing a global objective function, Ditto enables local devices to perform additional steps of SGD, within certain constraints, to minimize their individual losses. Ditto proves effective in enhancing testing accuracy among local devices and promoting fairness.

Our approach shares similarities with *FedMGDA* + Hu et al. (2022), which treats the FL task as a multi-objective optimization problem. The objective here is to simultaneously minimize the loss function of each FL client. To ensure that the performance of any client is not compromised, *FedMGDA*+
leverages Pareto-stationary solutions to identify a common descent direction for all selected clients. In a similar approach, Hamidi & YANG (2024); Mohajer Hamidi & Damen (2024) use ideas from multi-objective optimization to ensure fairness in FL models.