

## Appendix for *Semi-Variance Reduction for Fair Federated Learning*

### A PROOFS

**Lemma 1.** For any model parameter  $\theta$ , the gradient of the global objective  $F(\theta)$  defined in equation 3 can be expressed as

$$\nabla F(\theta) = \sum_{i=1}^n w_i(\theta) \nabla f_i(\theta), \quad w_i(\theta) = \frac{1}{n} + \frac{2\beta(f_i(\theta) - \bar{f}(\theta))}{n}, \quad \bar{f}(\theta) = \frac{\sum_i f_i(\theta)}{n}. \quad (8)$$

*Proof.* From equation 3 and with  $p_i = \frac{1}{n}$ , we have:

$$\begin{aligned} n\nabla F(\theta) &= \sum_i \nabla f_i(\theta) + 2\beta \sum_i \left[ (f_i(\theta) - \bar{f}(\theta)) (\nabla f_i(\theta) - \nabla \bar{f}(\theta)) \right] \\ &= \sum_i \nabla f_i(\theta) + 2\beta \sum_i \left[ (f_i(\theta) - \bar{f}(\theta)) \nabla f_i(\theta) - (f_i(\theta) - \bar{f}(\theta)) \nabla \bar{f}(\theta) \right] \\ &= \sum_i \left( 1 + 2\beta(f_i(\theta) - \bar{f}(\theta)) \right) \nabla f_i(\theta) - 2\beta \sum_i (f_i(\theta) - \bar{f}(\theta)) \nabla \bar{f}(\theta) \\ &= \sum_i \left( 1 + 2\beta(f_i(\theta) - \bar{f}(\theta)) \right) \nabla f_i(\theta) \end{aligned} \quad (15)$$

Hence,

$$\nabla F(\theta) = \sum_i \frac{1 + 2\beta(f_i(\theta) - \bar{f}(\theta))}{n} \nabla f_i(\theta) \quad (16)$$

□

### Derivation of equation 6

$$\begin{aligned} \sum_i f_i + \beta \sum_i \left| f_i(\theta) - \frac{1}{n} \sum_j f_j(\theta) \right|^2 &= \sum_i f_i + \beta \sum_i \left| \frac{n-1}{n} f_i(\theta) - \frac{1}{n} \sum_{j \neq i} f_j(\theta) \right|^2 \\ &= \sum_i f_i + \frac{\beta}{n^2} \sum_i \left| \sum_{j \neq i} (f_i(\theta) - f_j(\theta)) \right|^2 \\ &\leq \sum_i f_i + \frac{\beta}{n^2} \sum_i \sum_{j \neq i} \left| f_i(\theta) - f_j(\theta) \right|^2 \\ &= \sum_i f_i + \frac{2\beta}{n^2} \sum_{j \neq i} \left| f_i(\theta) - f_j(\theta) \right|^2. \end{aligned} \quad (17)$$

**Lemma 2.** In each communication round between the clients and the server, let  $>_C$  denote the set of clients whose local loss function is greater than the average loss function  $\bar{f}(\theta)$ . For any model parameter  $\theta$ , the gradient of the global objective  $F(\theta)$  defined in equation 9 can be expressed as

$$\nabla F(\theta) = \sum_{i=1}^n w_i(\theta) \nabla f_i(\theta), \quad (11)$$

where:

$$\bar{f}(\theta) = \frac{\sum_i f_i(\theta)}{n}, \quad w_i(\theta) = \begin{cases} \frac{1}{n} + \frac{2\beta(f_i(\theta) - \bar{f}(\theta))}{n} - \frac{2\beta \sum_{j \in >c} (f_j(\theta) - \bar{f}(\theta))}{n^2}, & \text{if } i \in >c \\ \frac{1}{n} - \frac{2\beta \sum_{j \in >c} (f_j(\theta) - \bar{f}(\theta))}{n^2}, & \text{if } i \notin >c \end{cases} \quad (12)$$

*Proof.* From equation 9 and with  $p_i = \frac{1}{n}$ , we have:

$$\begin{aligned} n\nabla F(\theta) &= \sum_i \nabla f_i(\theta) + 2\beta \sum_{i \in >c} \left[ (f_i(\theta) - \bar{f}(\theta)) (\nabla f_i(\theta) - \nabla \bar{f}(\theta)) \right] \\ &= \sum_i \nabla f_i(\theta) + 2\beta \sum_{i \in >c} \left[ (f_i(\theta) - \bar{f}(\theta)) \nabla f_i(\theta) - (f_i(\theta) - \bar{f}(\theta)) \nabla \bar{f}(\theta) \right] \\ &= \sum_{i \notin >c} \nabla f_i(\theta) + \sum_{i \in >c} \left( 1 + 2\beta(f_i(\theta) - \bar{f}(\theta)) \right) \nabla f_i(\theta) - 2\beta \sum_{i \in >c} (f_i(\theta) - \bar{f}(\theta)) \nabla \bar{f}(\theta) \end{aligned} \quad (18)$$

The last term in the above equation can be written as:

$$\begin{aligned} &- 2\beta \sum_{i \in >c} (f_i(\theta) - \bar{f}(\theta)) \nabla \bar{f}(\theta) \\ &= - \left[ \frac{2\beta}{n} \left( \sum_{i \in >c} f_i(\theta) \right) \times \left( \sum_j \nabla f_j(\theta) \right) \right] + \left[ \frac{2\beta}{n} \left( \sum_{i \in >c} \bar{f}(\theta) \right) \times \left( \sum_j \nabla f_j(\theta) \right) \right] \\ &= - \left[ \frac{2\beta}{n} \left( \sum_{i \in >c} f_i(\theta) - \bar{f}(\theta) \right) \times \left( \sum_j \nabla f_j(\theta) \right) \right] \end{aligned} \quad (19)$$

Hence,

$$\begin{aligned} n\nabla F(\theta) &= \sum_{i \in >c} \left( 1 + 2\beta(f_i(\theta) - \bar{f}(\theta)) - \frac{2\beta}{n} \left( \sum_{j \in >c} f_j(\theta) - \bar{f}(\theta) \right) \right) \nabla f_i(\theta) \\ &\quad + \sum_{i \notin >c} \left( 1 - \frac{2\beta}{n} \left( \sum_{j \in >c} f_j(\theta) - \bar{f}(\theta) \right) \right) \nabla f_i(\theta) \end{aligned} \quad (20)$$

Therefore,

$$\begin{aligned} \nabla F(\theta) &= \sum_{i \in >c} \left( \frac{1 + 2\beta(f_i(\theta) - \bar{f}(\theta)) - \frac{2\beta}{n} \left( \sum_{j \in >c} f_j(\theta) - \bar{f}(\theta) \right)}{n} \right) \nabla f_i(\theta) \\ &\quad + \sum_{i \notin >c} \left( \frac{1 - \frac{2\beta}{n} \left( \sum_{j \in >c} f_j(\theta) - \bar{f}(\theta) \right)}{n} \right) \nabla f_i(\theta) \end{aligned} \quad (21)$$

□

**Lemma 3.** Assuming  $P_i(x, y) = P_i(x|y)P_i(y) = P(x|y)P_i(y)$  for  $i \in \{1, \dots, n\}$ , for any model parameter  $\theta$ , Semi-VRRed global objective  $F(\theta)$  defined in equation 9 can be expressed as

$$F(\theta) = \sum_{j=1}^C \bar{P}(j) \bar{\ell}_j(\theta) + \frac{\beta}{n} \sum_{i=1}^n \left( \sum_{j=1}^C [P_i(j) - \bar{P}(j)] \bar{\ell}_j(\theta) \right)_+^2, \quad (13)$$

where  $\bar{P}(j) = \frac{\sum_{i=1}^n P_i(j)}{n}$  is the marginal distribution of class  $j$  in the global dataset.

*Proof.* From equation 9 and with  $p_i = \frac{1}{n}$ , we have:

$$\begin{aligned}
\bar{f}(\theta) &= \sum_{i=1}^n \frac{f_i(\theta)}{n} = \frac{1}{n} \sum_{i=1}^n \left[ \mathbb{E}_{(x,y) \sim p_i(x,y)} [\ell(h(x, \theta), y)] \right] \\
&= \frac{1}{n} \sum_{i=1}^n \left[ \sum_{j=1}^C p_i(j) \times \mathbb{E}_{(x,y) \sim p(x|y=j)} [\ell(h(x, \theta), j)] \right] \\
&= \frac{1}{n} \sum_{i=1}^n \left[ \sum_{j=1}^C p_i(j) \bar{\ell}_j(\theta) \right] \\
&= \sum_{j=1}^C \left[ \left( \frac{\sum_{i=1}^n p_i(j)}{n} \right) \bar{\ell}_j(\theta) \right] \\
&= \sum_{j=1}^C \bar{p}(j) \bar{\ell}_j(\theta), \tag{22}
\end{aligned}$$

where  $\bar{p}(j) = \frac{\sum_{i=1}^n p_i(j)}{n}$  is the ratio of data points with label  $j$  in the global dataset. Similarly, we have:

$$f_i(\theta) = \sum_{j=1}^C p_i(j) \bar{\ell}_j(\theta). \tag{23}$$

By plugging in the above equivalences for  $f_i(\theta)$  and  $\bar{f}(\theta)$  into equation 9, we get to equation 13.  $\square$

---

**Algorithm 2:** FedAvg

---

- 1 **Input:** global epoch  $T$ , client number  $n$ , loss function  $f_i$ , number of samples  $n_i$  for client  $i$ , number of total samples  $N$ , initial global model  $\theta_0$ , local step number  $K_i$  for client  $i$ , learning rate  $\eta$
  - 2 **for**  $t = 0, 1 \dots T - 1$  **do**
  - 3     randomly select  $\mathcal{S}_t \subseteq [n]$
  - 4      $\theta_t^{(i)} = \theta_t$  for  $i \in \mathcal{S}_t$ ,  $N = \sum_{i \in \mathcal{S}_t} n_i$
  - 5     **for**  $i$  in  $\mathcal{S}_t$  **do** // in parallel
  - 6         starting from  $\theta_t^{(i)}$ , take  $K_i$  local SGD steps on  $f_i$ , with learning rate  $\eta$ , to find  $\theta_{t+1}^{(i)}$
  - 7      $\theta_{t+1} = \sum_{i \in \mathcal{S}_t} \frac{n_i}{N} \theta_{t+1}^{(i)}$
  - 8 **Output:** global model  $\theta_T$
- 

In algorithm 2, we have reported the FedAvg algorithm for easier reference. Due to the similarity of VRed's objective function to that of FedAvg, we build its convergence proof on top of the proof for FedAvg in Zhang et al. (2022a). We borrow the following theorem on convergence of FedAvg from the same work. We refer the reader to the work for the detailed proof.

**Theorem 2 (FedAvg).** Denote  $p_i = \frac{n_i}{N}$ . Given Assumption 1, assume that the local learning rate satisfies  $\eta \leq \frac{1}{6LK_i}$  for any  $i \in [n]$  and

$$\eta \leq \frac{1}{L} \sqrt{\frac{1}{24(e-2)(\sum_i p_i^2)(\sum_i K_i^4)}}. \tag{24}$$

Running FedAvg for  $T$  global epochs we have:

$$\min_{0 \leq t \leq T-1} \mathbb{E} \|\nabla F(\theta_t)\|^2 \leq \frac{12}{(11\gamma - 9)\eta} \left( \frac{F_0 - F^*}{T} + \Psi_\sigma \right),$$

with  $\gamma = \sum_i p_i K_i$  for full participation and  $\gamma = \min_i K_i$  for partial participation,  $F_0 = F(\theta_0)$ ,  $F^* = \min_{\theta} F(\theta)$  the optimal value, and

$$\Psi_{\sigma} = \frac{\eta}{2} \left( \sum_{i=1}^n p_i^2 \right) \left[ \sum_{i=1}^n K_i^2 (\sigma_{l,i}^2 + 2\sigma_g^2) + 2(e-2)\eta^2 L^2 \sum_{i=1}^n K_i^3 (\sigma_{l,i}^2 + 6K_i \sigma_g^2) \right].$$

Based on the above theorem for FedAvg, we now prove the convergence of our VRed algorithm.

**Theorem 1 (VRed).** Denote  $\tilde{L} = (L + \beta ML + 2\beta L_0^2)$  and  $p_i = \frac{n_i}{N}$ . Given Assumptions 1 and 2, assume that the local learning rate satisfies  $\eta \leq \frac{1}{6\tilde{L}K_i}$  for any  $i \in [n]$  and:

$$\eta \leq \frac{1}{L} \sqrt{\frac{1}{24(e-2)(\sum_i p_i^2)(\sum_i K_i^4)}} \quad (14)$$

By running Algorithm 1 for  $T$  global epochs we have:

$$\min_{0 \leq t \leq T-1} \mathbb{E} \|\nabla F(\theta_t)\|^2 \leq \frac{12}{(11\gamma - 9)\eta} \left( \frac{F_0 - F^*}{T} + \tilde{\Psi}_{\sigma} \right),$$

with  $\gamma = \sum_i p_i K_i$  for full participation and  $\gamma = \min_i K_i$  for partial participation,  $F_0 = F(\theta_0)$ ,  $F^* = \min_{\theta} F(\theta)$  the optimal value, and

$$\tilde{\Psi}_{\sigma} = \frac{\eta}{2} \left( \sum_{i=1}^n p_i^2 \right) \left[ \sum_{i=1}^n K_i^2 (\tilde{\sigma}_{l,i}^2 + 2\tilde{\sigma}_g^2) + 2(e-2)\eta^2 L^2 \sum_{i=1}^n K_i^3 (\tilde{\sigma}_{l,i}^2 + 6K_i \tilde{\sigma}_g^2) \right],$$

where  $\tilde{\sigma}_{l,i}^2 = (2(2\beta\mu + 1)^2 \sigma_{l,i}^2 + 8\beta^2 M^2 L_0^2)$  and  $\tilde{\sigma}_g^2 = ((2\beta\mu + 1)\sigma_g^2 + 2\beta M L_0)$  and  $\mu = \sum_i p_i f_i(\theta_0)$ .

*Proof.* We first rewrite a simplified version of the VRed objective function (equation 3) in the following.

$$F(\theta) = \sum_i p_i G_i(\theta) = \sum_i p_i (f_i(\theta) + \beta (f_i(\theta) - \mu)^2), \quad (25)$$

where  $G_i(\theta) = f_i(\theta) + \beta (f_i(\theta) - \mu)^2$ , and also,  $\mu = \sum_i p_i f_i$  is fixed during clients local computations. In the beginning of each communication round, we update  $\mu$  for the next round of local computations. In other words:

$$\mu_{t+1} = \sum_{i=1}^n p_i f_i(\theta_{t+1}). \quad (26)$$

With these notations, it suffices to find the constants in Assumption 1 for  $G_i(\theta)$ .

**Smoothness** We have:

$$\begin{aligned} \nabla^2 G_i(\theta) &= \nabla^2 f_i(\theta) + 2\beta((f_i(\theta) - \mu)\nabla^2 f_i(\theta) + \nabla f_i(\theta)\nabla f_i(\theta)^{\top}) \\ &\preceq L + 2\beta\left(\frac{M}{2}\nabla^2 f_i(\theta) + \nabla f_i(\theta)\nabla f_i(\theta)^{\top}\right) \\ &\preceq (L + \beta ML + 2\beta L_0^2)I, \end{aligned} \quad (27)$$

where in the last line we used Assumption 2 and the following:

$$\begin{aligned} \|\nabla f_i(\theta)\nabla f_i(\theta)^{\top}\|_{sp} &= \sup_{\|u\|=1} \sup_{\|v\|=1} \langle \nabla f_i(\theta)\nabla f_i(\theta)^{\top} u; v \rangle \\ &= \sup_{\|u\|=1} \sup_{\|v\|=1} (\nabla f_i(\theta)^{\top} u)^{\top} \nabla f_i(\theta)^{\top} v = \|\nabla f_i(\theta)\|^2 \leq L_0^2, \end{aligned} \quad (28)$$

where in the second line we used Cauchy-Schwarz inequality and Assumption 2. Hence, from eq. (27), we conclude that  $G_i(\theta)$  is Lipschitz smooth:

$$\|\nabla G_i(\theta) - \nabla G_i(\theta')\| \leq (L + \beta ML + 2\beta L_0^2)\|\theta - \theta'\|. \quad (29)$$

**Variance constants** For the global variance between clients gradients, we have:

$$\begin{aligned}
& \|\nabla G_i(\theta) - \nabla G_j(\theta)\| \\
&= \left\| \left( \nabla f_i(\theta) + 2\beta(f_i(\theta) - \mu)\nabla f_i(\theta) \right) - \left( \nabla f_j(\theta) + 2\beta(f_j(\theta) - \mu)\nabla f_j(\theta) \right) \right\| \\
&\leq \|\nabla f_i(\theta) - \nabla f_j(\theta)\| + 2\beta\|(f_i(\theta) - \mu)\nabla f_i(\theta) - (f_j(\theta) - \mu)\nabla f_j(\theta)\| \\
&\leq \|\nabla f_i(\theta) - \nabla f_j(\theta)\| + 2\beta\|f_i(\theta)\nabla f_i(\theta) - f_j(\theta)\nabla f_j(\theta)\| + 2\beta\mu\|\nabla f_i(\theta) - \nabla f_j(\theta)\| \\
&\leq (2\beta\mu + 1)\|\nabla f_i(\theta) - \nabla f_j(\theta)\| + 2\beta\|f_i(\theta)\nabla f_i(\theta) - f_j(\theta)\nabla f_j(\theta)\| \\
&\leq (2\beta\mu + 1)\|\nabla f_i(\theta) - \nabla f_j(\theta)\| + 2\beta f_i(\theta)\|\nabla f_i(\theta)\| + 2\beta f_j(\theta)\|\nabla f_j(\theta)\| \\
&\leq (2\beta\mu + 1)\sigma_g + 2\beta f_i(\theta)\|\nabla f_i(\theta)\| + 2\beta f_j(\theta)\|\nabla f_j(\theta)\| \\
&\leq (2\beta\mu + 1)\sigma_g + \beta M\|\nabla f_i(\theta)\| + \beta M\|\nabla f_j(\theta)\| \\
&\leq (2\beta\mu + 1)\sigma_g + 2\beta M L_0,
\end{aligned} \tag{30}$$

where in line seven we used Assumption 1, and in lines eight and nine, we used Assumption 2. Therefore, we have:

$$\|\nabla G_i(\theta) - \nabla G_j(\theta)\|^2 \leq \left( (2\beta\mu + 1)\sigma_g + 2\beta M L_0 \right)^2. \tag{31}$$

For the local variance term, we define  $\varphi(t) = t + \beta(t - \mu)^2$ . Similar to the derivation of equation 30, we have:

$$\begin{aligned}
& \|\nabla G_i(\theta) - \nabla(\varphi \circ \ell_S)(\theta)\| \\
&= \left\| \left( \nabla f_i(\theta) + 2\beta(f_i(\theta) - \mu)\nabla f_i(\theta) \right) - \left( \nabla \ell_S(\theta) + 2\beta(\ell_S(\theta) - \mu)\nabla \ell_S(\theta) \right) \right\| \\
&\leq (2\beta\mu + 1)\|\nabla f_i(\theta) - \nabla \ell_S(\theta)\| + 2\beta f_i(\theta)\|\nabla f_i(\theta)\| + 2\beta \ell_S(\theta)\|\nabla \ell_S(\theta)\| \\
&\leq (2\beta\mu + 1)\|\nabla f_i(\theta) - \nabla \ell_S(\theta)\| + 2\beta M L_0,
\end{aligned} \tag{32}$$

where in line four, we used Assumption 2. By taking the square on both sides and the expectation over  $S \sim \mathcal{B}_i^b$ , we get:

$$\begin{aligned}
\mathbb{E}_{S \sim \mathcal{B}_i^b} \|\nabla G_i(\theta) - \nabla(\varphi \circ \ell_S)(\theta)\|^2 &\leq \mathbb{E}_{S \sim \mathcal{B}_i^b} \left( (2\beta\mu + 1)\|\nabla f_i(\theta) - \nabla \ell_S(\theta)\| + 2\beta M L_0 \right)^2 \\
&\leq \mathbb{E}_{S \sim \mathcal{B}_i^b} \left( 2(2\beta\mu + 1)^2 \|\nabla f_i(\theta) - \nabla \ell_S(\theta)\|^2 + 8\beta^2 M^2 L_0^2 \right) \\
&= \left( 2(2\beta\mu + 1)^2 \sigma_{l,i}^2 + 8\beta^2 M^2 L_0^2 \right).
\end{aligned} \tag{33}$$

In the third line, we used  $(a + b)^2 \leq 2(a^2 + b^2)$ . We also used Assumption 1 in the same line.  $\square$

## B EXAMPLES

We borrow the following example on class imbalance in FL from Shen et al. (2022) to provide a better understanding of lemma 3. The following example shows an extreme class imbalance, which Semi-VRed can handle efficiently.

**Example 1.** Let  $u$  be the uniform distribution over the existing  $C$  classes. Also, let  $\delta_c$  be the Dirac distribution of class  $c$ . Now, without loss of generality, let's assume that  $C = 2$  (binary classification problem). For the  $n$  existing clients, we have:

$$p_i(y) = \begin{cases} \alpha u + (1 - \alpha)\delta_1 & \text{if } i = 1 \\ \alpha u + (1 - \alpha)\delta_2 & \text{if } i \in \{2, \dots, n\} \end{cases} \tag{34}$$

Accordingly, we have:

$$p_i(1) = \begin{cases} 1 - \frac{\alpha}{2} & \text{if } i = 1 \\ \frac{\alpha}{2} & \text{if } i \in \{2, \dots, n\} \end{cases} \quad (35)$$

$$p_i(2) = \begin{cases} \frac{\alpha}{2} & \text{if } i = 1 \\ 1 - \frac{\alpha}{2} & \text{if } i \in \{2, \dots, n\} \end{cases} \quad (36)$$

Therefore,

$$f_i(\theta) = \begin{cases} (1 - \frac{\alpha}{2})\bar{\ell}_1(\theta) + \frac{\alpha}{2}\bar{\ell}_2(\theta) & \text{if } i = 1 \\ \frac{\alpha}{2}\bar{\ell}_1(\theta) + (1 - \frac{\alpha}{2})\bar{\ell}_2(\theta) & \text{if } i \in \{2, \dots, n\} \end{cases} \quad (37)$$

Hence,

$$\bar{f}(\theta) = \left(\frac{\alpha}{2} + \frac{1-\alpha}{n}\right)\bar{\ell}_1(\theta) + \left(\frac{\alpha}{2} + \frac{(1-\alpha)(n-1)}{n}\right)\bar{\ell}_2(\theta) \quad (38)$$

Clearly, we can see that if  $\alpha \approx 0$  and  $n$  is large, then  $\bar{\ell}_1(\theta)$ , which is the loss over the minority data will have a small weight, which leads to  $\bar{\ell}_1(\theta)$  being larger than  $\bar{\ell}_2(\theta)$  and poor performance on the minority class 1. Now, if we rewrite the *Semi-VRed* objective function (equation 9), we have:

$$F(\theta) = \left(\frac{\alpha}{2} + \frac{1-\alpha}{n}\right)\bar{\ell}_1(\theta) + \left(\frac{\alpha}{2} + \frac{(1-\alpha)(n-1)}{n}\right)\bar{\ell}_2(\theta) + \frac{\beta(n-1)^2(1-\alpha)^2}{n^3} \left(\bar{\ell}_1(\theta) - \bar{\ell}_2(\theta)\right)^2 \quad (39)$$

For  $\alpha \approx 0$ :

$$F(\theta) \approx \bar{\ell}_2(\theta) + \frac{\beta}{n} \left(\bar{\ell}_1(\theta) - \bar{\ell}_2(\theta)\right)^2, \quad (40)$$

which improves  $\bar{\ell}_2(\theta)$ , thanks to its regularization term. Hence, the performance of client 1 and consequently, fairness in the system will improve.

## C EXPERIMENTAL SETUP

In this section, we provide more experimental details that are deferred from the main paper. For each experiment, we report the average result obtained from three runs with different random seeds. For our experiments, we used an internal GPU server with six NVIDIA Tesla P100. The experiments last about 4 weeks in total.

### C.1 DATASETS AND MODELS

In this subsection, we describe the datasets we use in our experiments. For all the datasets we use a batch size of 64.

**CIFAR-10/100** (Krizhevsky et al., 2009) are two image classification datasets vastly used in the literature as benchmark datasets. Each of these datasets contains 50000 sample images with 10/100 balanced classes for CIFAR-10 and CIFAR-100, respectively. We use Dirichlet allocation (Wang et al., 2019) to distribute the data among 50 clients with label shift: we split the set of samples from class  $k$  ( $\mathcal{S}_k$ ) to  $n$  subsets ( $\mathcal{S}_k = \mathcal{S}_{k,1} \cup \mathcal{S}_{k,2} \cup \dots \cup \mathcal{S}_{k,n}$ ), where  $n$  is the number of clients and  $\mathcal{S}_{k,j}$  corresponds to the client  $j$ . We do the split based on Dirichlet distribution with parameter 0.05 ( $\text{Dir}(0.05)$ ). When the split is done for all classes, we gather the samples corresponding to each client from different classes: assuming there are  $C$  classes in total  $\mathcal{S}_{1,j} \cup \mathcal{S}_{2,j} \cup \dots \cup \mathcal{S}_{C,j}$  is the data allocated to the client  $j$ . Having allocated the data of each client, we split them into train and

test set for each client. The train-test split ratio is 50-50 and 60-40 for CIFAR-10 and CIFAR-100, respectively.

**CINIC-10** (Darlow et al., 2018) is another benchmark vision dataset that we use in our experiments. There are a total of 270,000 sample images, which we distribute with label shift between 50 clients based on  $\text{Dir}(0.5)$  distribution Wang et al. (2019). We then randomly split the data of each client into train and test sets with split ratio 50-50.

**StackOverflow** (The Tensorflow Federated Authors, 2019) is a language dataset consisting of Shakespeare dialogues for the task of next word prediction. There is a natural heterogeneous partition of the dataset and we treat each speaking role as a client. We filter out the clients (speaking roles) with less than 10,000 samples from the original dataset and randomly select 20 clients from the remaining. Finally, we split the data of each client into train and test sets with a ratio of 50-50.

Table 2 provides a summary of the datasets we used and the models used for each of them.

Table 2: Details of the experiments and the datasets. ResNet-18: residual neural network (He et al., 2016). GN: Group Normalization (Wu & He, 2018); RNN: Recurrent Neural Network; LSTM: Long Short-Term Memory layer; FC: fully connected layer.

Dataset	Train samples	Test samples	Partition method	clients	Model
CIFAR-10	24959	25041	$\text{Dir}(0.05)$	50	ResNet-18 + GN
CIFAR-100	39445	10555	$\text{Dir}(0.05)$	50	ResNet-18 + GN
CINIC-10	134713	134966	$\text{Dir}(0.5)$	50	ResNet-18 + GN
StackOverflow	109671	109621	realistic partition	20	RNN (1 LSTM + 2 FC)

## C.2 ALGORITHMS AND THEIR HYPERPARAMETERS

We use most recent fair FL algorithms existing in the literature as our baseline algorithms, including: FedAvg (McMahan et al., 2017),  $q$ -FFL (Li et al., 2020c), AFL (Mohri et al., 2019), PropFair (Zhang et al., 2022a), TERM (Li et al., 2020a), GiFair (Yue et al., 2021). For each pair of dataset and algorithm, we find the best learning rate based on a grid search. In the following, we have reported the learning rate grid we use for each dataset:

- CIFAR-10:  $\{1e-3, 2e-3, 5e-3, 1e-2, 2e-2, 5e-2\}$ ;
- CIFAR-100:  $\{1e-3, 2e-3, 5e-3, 1e-2, 2e-2, 5e-2\}$ ;
- CINIC-10:  $\{1e-3, 2e-3, 5e-3, 1e-2, 2e-2, 5e-2\}$ ;
- StackOverflow:  $\{1e-2, 5e-2, 1e-1, 5e-1, 1\}$ .

The best learning rate used for each (dataset, algorithm) pair is reported in Table 4.

We now explain the algorithms we use and how we tune their hyperparameters. We adapt TERM with only client-level fairness ( $\alpha > 0$ ) and no sample-level fairness ( $\tau = 0$ ). We tune the hyperparameter  $\alpha$  for each dataset based on a grid search in the grid  $\{0.01, 0.1, 0.5, 1\}$ . We have reported the best value of  $\alpha$  for each dataset in Table 5. For AFL, there are two hyperparameters:  $\gamma_w$  and  $\gamma_\lambda$ . We tune the learning rate  $\gamma_w$  from the corresponding grid and choose the default value  $\gamma_\lambda = 0.1$ . For  $q$ -FFL, we use the  $q$ -FedAvg algorithm (Li et al., 2020c). We also tune the hyperparameter  $q$  from the grid  $\{0.01, 0.1, 1\}$ . We find that for all the used datasets,  $q = 0.1$  has the best performance (as reported in Table 5). We also tried larger values out of the grid and they often lead to divergence of the  $q$ -FFL algorithm. We adopt the Global GiFair model (Yue et al., 2021), which results in a single global model. In order to have client-level fairness, we treat each client as a group of size 1. For tuning the regularization weight of GiFair ( $\lambda$ ), we follow (Yue et al., 2021). As stated in the paper, there is an upper-bound for  $\lambda$  (see §3 in the paper). For our experiments, the upper-bound is  $\lambda \leq \min_i \{\frac{w_i}{n-1}\}$ , where  $w_i$  is the ratio of total samples allocated to the client  $i$  and  $n$  is the number of clients. We try four different values in the interval and choose the best one. When the number of clients is large, this upper-bound is small, and for all of our datasets, this upper-bound was the best value, as reported in Table 5. We tune  $M$  for the PropFair algorithm based on a grid search in  $\{2, 3, 4, 5\}$ . Finally, for our VRed and Semi-VRed algorithms, we tune the regularization weight  $\beta$  based on grid search on the grid  $\{0.01, 0.05, 0.1, 0.2, 0.5, 1\}$ . Larger values of  $\beta$  often resulted in

Table 3: Details of the existing fairfl algorithms.  $f_i$  is the loss function of the client  $i$ .

FL algorithm	Objective	Reference
FedAvg	$\sum_i f_i$	McMahan et al. (2017)
AFL	$\max_i f_i$	Mohri et al. (2019)
q-FFL	$\sum_i f_i^{q+1}$	Li et al. (2020c)
TERM	$\sum_i e^{\alpha f_i}$	Li et al. (2020a)
PropFair	$-\sum_i \log(M - f_i)$	Zhang et al. (2022a)
GiFair	$\sum_i f_i + \lambda \sum_{i \neq j}  f_i - f_j $	Yue et al. (2021)
VRed	$\sum_i f_i + \beta \sum_i \left( f_i(\theta) - \frac{1}{n} \sum_j f_j(\theta) \right)^2$	this work
Semi-VRed	$\sum_i f_i + \beta \sum_i \left( f_i(\theta) - \frac{1}{n} \sum_j f_j(\theta) \right)_+^2$	this work

Table 4: The best learning rates used for training each algorithm on different datasets.

Datasets	FedAvg	q-FFL	AFL	TERM	PropFair	GiFair	VRed	Semi-VRed
CIFAR-10	5e-3	5e-3	5e-3	1e-2	1e-2	5e-3	5e-3	5e-3
CIFAR-100	2e-3	2e-3	5e-3	1e-2	1e-2	5e-3	5e-3	5e-3
CINIC-10	1e-2	5e-3	1e-2	1e-2	2e-2	2e-2	5e-3	5e-3
StackOverflow	2e-1	5e-2	5e-2	2e-1	5e-1	2e-1	5e-1	5e-1

the divergence of the algorithms. We have reported the best value of all of the hyperparameters for each dataset in Table 5.

Table 5: The best values of hyperparameters used for different datasets, chosen based on grid search.

Algorithm	CIFAR-10	CIFAR-100	CINIC-10	StackOverflow
<b>q-FFL</b> $q$	1e-1	1e-1	1e-1	1e-1
<b>TERM</b> $\alpha$	1e-2	5e-1	5e-1	5e-1
<b>GiFair</b> $\lambda$	6e-5	2.6e-4	5e-5	2.4e-3
<b>PropFair</b> $M$	3	3	5	4
<b>VRed</b> $\beta$	5e-1	1e-1	2e-1	1e-1
<b>Semi-VRed</b> $\beta$	5e-1	1e-2	2e-1	2e-1

### C.3 DETAILED RESULTS

In Table 6, we report detailed results obtained from the algorithms we study in this work. We use a default batch size of 64 for all the experiments. The statistics we report include: 1. the average test accuracy across clients (overall average performance) 2. the standard deviation of test accuracies across clients 3. the lowest (worst) test accuracy among clients 4. the lowest 10% test accuracies 5. the lowest 20% test accuracies 6. the highest 10% test accuracies. For each experiment, we report the average result obtained from three runs with different random seeds. As can be observed, our proposed algorithms VRed and Semi-VRed consistently beat almost all the baseline algorithms across different datasets in terms of various fairness metrics. Also, Semi-VRed can improve the overall average performance (mean test accuracy) for three of the datasets as well.

Following Figure 1, we have compared our proposed algorithms with the baseline algorithms in terms of their worst 20% test accuracies as well and the visualized results are shown in Figure 2.



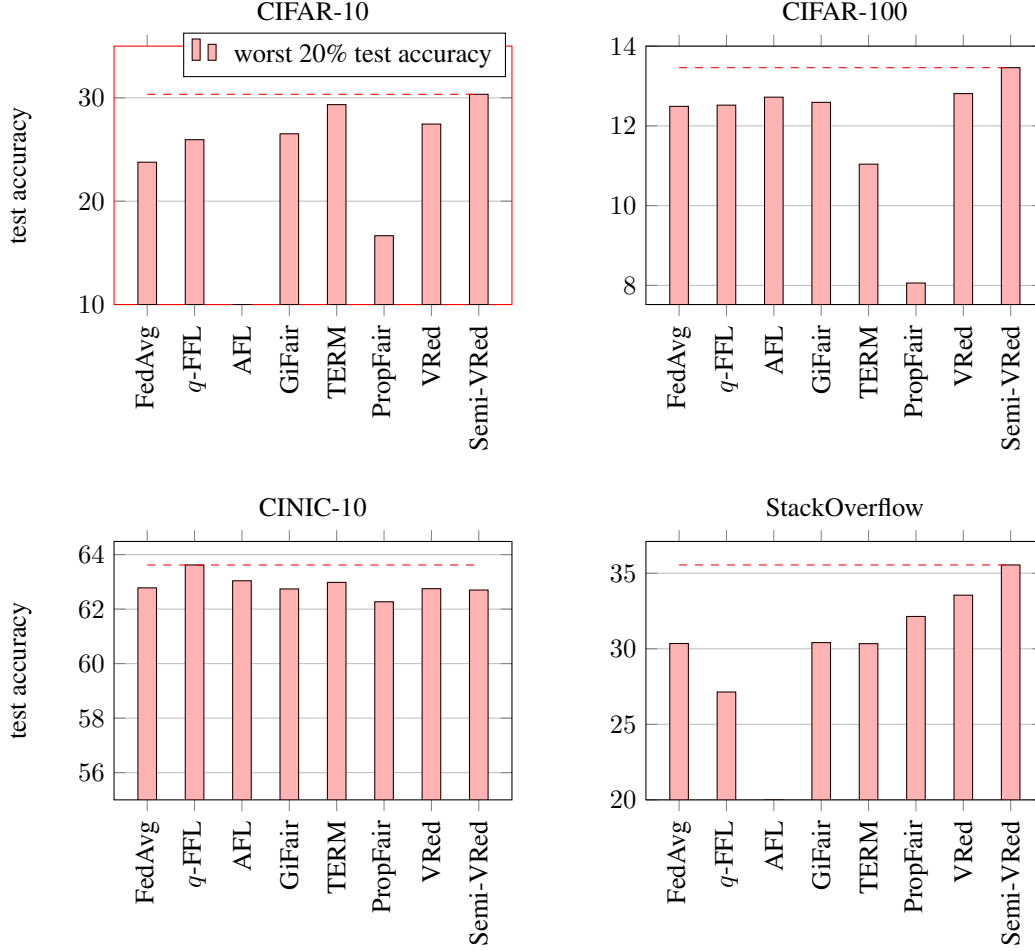


Figure 2: Worst 20% test accuracies for different algorithms. **top left:** CIFAR-10, **top right:** CIFAR-100, **bottom left:** CINIC-10, **bottom right:** StackOverflow. Due to divergence, results for AFL on CIFAR-10 and StackOverFlow are not shown. All subfigures share the same legends and axis labels.

Table 6: Comparison among federated learning algorithms on CIFAR-10, CIFAR-100, CINIC-10 and StackOverflow datasets with test accuracies (%) from clients. All algorithms are fine-tuned. **Mean**: the average test accuracy across all clients; **Std**: standard deviation of clients test accuracies; **Worst**: the worst test accuracy among clients; **Worst (10/20%)**: the worst 10/20% test accuracies of clients; **Best (10%)**: the best 10% test accuracies of clients.

Dataset	Algorithm	Mean	Std	Worst	Worst (10%)	Worst (20%)	Best (10%)
CIFAR-10	FedAvg, Ditto	43.45 $\pm$ 0.60	14.33 $\pm$ 0.62	9.35 $\pm$ 3.13	18.86 $\pm$ 0.99	23.77 $\pm$ 0.70	68.97 $\pm$ 0.81
	<i>q</i> -FFL	45.46 $\pm$ 0.74	14.31 $\pm$ 2.03	18.71 $\pm$ 3.36	21.23 $\pm$ 3.06	25.95 $\pm$ 3.51	72.31 $\pm$ 2.88
	AFL	-	-	-	-	-	-
	GiFair	45.05 $\pm$ 0.64	12.93 $\pm$ 0.44	16.79 $\pm$ 3.55	22.65 $\pm$ 2.03	26.52 $\pm$ 0.76	65.62 $\pm$ 2.59
	TERM	<b>45.61</b> $\pm$ 1.03	<b>12.24</b> $\pm$ 0.56	13.80 $\pm$ 5.25	24.89 $\pm$ 1.37	29.34 $\pm$ 0.61	68.65 $\pm$ 1.27
	PropFair	36.95 $\pm$ 0.21	15.16 $\pm$ 1.33	1.14 $\pm$ 1.62	12.49 $\pm$ 0.28	16.66 $\pm$ 1.31	66.04 $\pm$ 4.24
	VRed	44.43 $\pm$ 0.88	13.05 $\pm$ 1.32	18.61 $\pm$ 3.12	24.28 $\pm$ 2.22	27.46 $\pm$ 1.56	69.31 $\pm$ 3.48
	Semi-VRed	45.47 $\pm$ 0.10	12.58 $\pm$ 0.23	<b>19.04</b> $\pm$ 6.73	<b>27.08</b> $\pm$ 1.76	<b>30.34</b> $\pm$ 1.05	<b>72.50</b> $\pm$ 0.88
CIFAR-100	FedAvg, Ditto	20.20 $\pm$ 0.31	6.50 $\pm$ 0.21	<b>10.36</b> $\pm$ 0.69	11.07 $\pm$ 0.54	12.49 $\pm$ 0.51	33.88 $\pm$ 0.09
	<i>q</i> -FFL	20.25 $\pm$ 0.11	6.30 $\pm$ 0.27	9.66 $\pm$ 0.33	11.09 $\pm$ 0.67	12.52 $\pm$ 0.46	33.96 $\pm$ 0.90
	AFL	18.98 $\pm$ 0.71	<b>4.91</b> $\pm$ 0.37	9.81 $\pm$ 0.69	11.31 $\pm$ 0.18	12.72 $\pm$ 0.21	28.68 $\pm$ 1.71
	GiFair	19.81 $\pm$ 0.32	5.74 $\pm$ 0.16	9.35 $\pm$ 0.34	11.19 $\pm$ 0.24	12.59 $\pm$ 0.49	32.30 $\pm$ 0.32
	TERM	18.00 $\pm$ 0.41	6.05 $\pm$ 0.18	8.86 $\pm$ 0.50	10.02 $\pm$ 0.44	11.04 $\pm$ 0.51	31.58 $\pm$ 0.98
	PropFair	14.97 $\pm$ 0.68	6.44 $\pm$ 0.34	5.40 $\pm$ 1.28	7.00 $\pm$ 1.11	8.06 $\pm$ 1.07	28.89 $\pm$ 0.91
	VRed	20.42 $\pm$ 0.36	6.08 $\pm$ 0.05	9.43 $\pm$ 1.01	11.21 $\pm$ 0.74	12.81 $\pm$ 0.85	33.59 $\pm$ 1.11
	Semi-VRed	<b>20.85</b> $\pm$ 0.39	6.26 $\pm$ 0.18	9.12 $\pm$ 1.47	<b>11.86</b> $\pm$ 0.74	<b>13.46</b> $\pm$ 0.63	<b>34.57</b> $\pm$ 1.20
CINIC-10	FedAvg, Ditto	86.26 $\pm$ 0.03	15.20 $\pm$ 0.07	50.48 $\pm$ 0.29	56.87 $\pm$ 0.36	62.78 $\pm$ 0.16	100.0 $\pm$ 0.00
	<i>q</i> -FFL	<b>86.63</b> $\pm$ 0.06	<b>14.88</b> $\pm$ 0.08	51.57 $\pm$ 0.82	57.77 $\pm$ 0.36	<b>63.62</b> $\pm$ 0.18	<b>100.0</b> $\pm$ 0.01
	AFL	86.45 $\pm$ 0.12	15.10 $\pm$ 0.11	51.57 $\pm$ 0.45	57.58 $\pm$ 0.29	63.04 $\pm$ 0.28	100.0 $\pm$ 0.00
	GiFair	86.28 $\pm$ 0.11	15.20 $\pm$ 0.13	49.66 $\pm$ 1.21	56.97 $\pm$ 0.29	62.74 $\pm$ 0.36	100.0 $\pm$ 0.00
	TERM	86.34 $\pm$ 0.04	15.12 $\pm$ 0.01	49.90 $\pm$ 0.42	57.21 $\pm$ 0.11	62.98 $\pm$ 0.04	100.0 $\pm$ 0.00
	PropFair	86.01 $\pm$ 0.17	15.34 $\pm$ 0.19	49.97 $\pm$ 1.23	56.53 $\pm$ 0.65	62.27 $\pm$ 0.55	99.99 $\pm$ 0.01
	VRed	85.79 $\pm$ 0.35	15.02 $\pm$ 0.06	51.57 $\pm$ 0.50	57.66 $\pm$ 0.30	62.75 $\pm$ 0.36	99.98 $\pm$ 0.01
	Semi-VRed	85.83 $\pm$ 0.33	14.95 $\pm$ 0.07	<b>51.59</b> $\pm$ 0.98	<b>58.00</b> $\pm$ 0.21	62.70 $\pm$ 0.14	99.96 $\pm$ 0.01
StackOverflow	FedAvg, Ditto	40.34 $\pm$ 0.06	6.98 $\pm$ 0.03	25.64 $\pm$ 0.11	27.12 $\pm$ 0.06	30.35 $\pm$ 0.03	49.70 $\pm$ 0.07
	<i>q</i> -FFL	37.79 $\pm$ 0.80	7.38 $\pm$ 0.09	22.54 $\pm$ 1.03	24.12 $\pm$ 1.00	27.14 $\pm$ 0.92	47.06 $\pm$ 0.66
	AFL	-	-	-	-	-	-
	TERM	40.34 $\pm$ 0.05	6.96 $\pm$ 0.06	25.56 $\pm$ 0.21	27.12 $\pm$ 0.20	30.41 $\pm$ 0.12	49.76 $\pm$ 0.10
	GiFair	40.34 $\pm$ 0.04	6.97 $\pm$ 0.03	25.71 $\pm$ 0.13	27.10 $\pm$ 0.11	30.34 $\pm$ 0.08	49.71 $\pm$ 0.09
	PropFair	41.76 $\pm$ 0.01	6.80 $\pm$ 0.05	27.30 $\pm$ 0.21	28.75 $\pm$ 0.19	32.14 $\pm$ 0.10	50.76 $\pm$ 0.08
	VRed	42.90 $\pm$ 0.05	6.64 $\pm$ 0.01	29.08 $\pm$ 0.09	<b>30.39</b> $\pm$ 0.05	33.55 $\pm$ 0.05	51.66 $\pm$ 0.03
	Semi-VRed	<b>42.90</b> $\pm$ 0.03	<b>6.60</b> $\pm$ 0.01	<b>29.10</b> $\pm$ 0.06	30.34 $\pm$ 0.09	<b>35.55</b> $\pm$ 0.05	<b>51.70</b> $\pm$ 0.04