

Inspired by the perturbed iterate analysis framework in [1], we first define the following auxiliary sequences for all  $t \geq 0$ :

1) If  $t = 0$ ,  $\tilde{\mathbf{w}}_m^{(t)} = \hat{\mathbf{w}}_m^{(0)}$ ; If  $t \geq 1$ ,  $\tilde{\mathbf{w}}_m^{(t)} = \tilde{\mathbf{w}}_m^{(t-1)} - \eta \nabla f_m(\hat{\mathbf{w}}_m^{(t-1)}; \mathcal{D}_m^{(t-1)})$ .

2)  $\mathbf{q}^{(t)} \triangleq \frac{1}{M} \sum_{m=1}^M \nabla f_m(\hat{\mathbf{w}}_m^{(t)}; \mathcal{D}_m^{(t)})$

3)  $\bar{\mathbf{q}}^{(t)} \triangleq \mathbb{E}_{\mathcal{D}_m^{(t)}}[\mathbf{q}^{(t)}] = \frac{1}{M} \sum_{m=1}^M \nabla f_m(\hat{\mathbf{w}}_m^{(t)})$

4)  $\tilde{\mathbf{w}}^{(t)} \triangleq \frac{1}{M} \sum_{m=1}^M \nabla f_m(\hat{\mathbf{w}}_m^{(t)}; \mathcal{D}_m^{(t)}) = \tilde{\mathbf{w}}^{(t-1)} - \eta \mathbf{q}^{(t-1)}$

5)  $\hat{\mathbf{w}}^{(t)} = \frac{1}{M} \sum_{m=1}^M \hat{\mathbf{w}}_m^{(t)}$

By the smoothness of  $F: \mathbb{R}^d \rightarrow \mathbb{R}$ , we have

$$\begin{aligned}
 & F(\tilde{\mathbf{w}}^{(t+1)}) - F(\tilde{\mathbf{w}}^{(t)}) \\
 & \leq \langle \nabla F(\tilde{\mathbf{w}}^{(t)}), \tilde{\mathbf{w}}^{(t+1)} - \tilde{\mathbf{w}}^{(t)} \rangle + \frac{L}{2} \|\tilde{\mathbf{w}}^{(t+1)} - \tilde{\mathbf{w}}^{(t)}\|^2 \\
 & = -\eta \langle \nabla F(\tilde{\mathbf{w}}^{(t)}), \mathbf{q}^{(t)} \rangle + \frac{\eta^2 L}{2} \|\mathbf{q}^{(t)}\|^2 \\
 & \stackrel{(a)}{\leq} -\eta \langle \nabla F(\tilde{\mathbf{w}}^{(t)}), \mathbf{q}^{(t)} \rangle + \eta^2 L \|\mathbf{q}^{(t)} - \bar{\mathbf{q}}^{(t)}\|^2 + \eta^2 L \|\bar{\mathbf{q}}^{(t)}\|^2 \\
 & = -\frac{\eta}{M} \sum_{m=1}^M \langle \nabla F(\tilde{\mathbf{w}}^{(t)}), \nabla f_m(\hat{\mathbf{w}}_m^{(t)}; \mathcal{D}_m^{(t)}) \rangle + \eta^2 L \|\mathbf{q}^{(t)} - \bar{\mathbf{q}}^{(t)}\|^2 + \eta^2 L \left\| \frac{1}{M} \sum_{m=1}^M \nabla f_m(\hat{\mathbf{w}}_m^{(t)}; \mathcal{D}_m^{(t)}) \right\|^2,
 \end{aligned} \tag{1}$$

where (a) is by Jensen's inequality. Taking expectation with respect to the sampling mini-batch  $\mathcal{D}_m^{(t)}$  by each edge device at time  $t$  gives

$$\begin{aligned}
 & \mathbb{E}[F(\tilde{\mathbf{w}}^{(t+1)})] - F(\tilde{\mathbf{w}}^{(t)}) \\
 & \stackrel{(a)}{\leq} -\frac{\eta}{2} (\|\nabla F(\tilde{\mathbf{w}}^{(t)})\|^2 + \left\| \frac{1}{M} \sum_{m=1}^M \nabla f_m(\hat{\mathbf{w}}_m^{(t)}) \right\|^2) + \frac{\eta}{2} \|\nabla F(\tilde{\mathbf{w}}^{(t)}) - \frac{1}{M} \sum_{m=1}^M \nabla f_m(\hat{\mathbf{w}}_m^{(t)})\|^2 + \eta^2 L \left\| \frac{1}{M} \sum_{m=1}^M \nabla f_m(\hat{\mathbf{w}}_m^{(t)}) \right\|^2 + \frac{\eta^2 L \sigma^2}{Mb^{(t)}} \\
 & \stackrel{(b)}{\leq} -\frac{\eta}{2M} \sum_{m=1}^M (\|\nabla F(\tilde{\mathbf{w}}^{(t)})\|^2 - L^2 \|\tilde{\mathbf{w}}^{(t)} - \hat{\mathbf{w}}_m^{(t)}\|^2) + \frac{2\eta^2 L - \eta}{2} \left\| \frac{1}{M} \sum_{m=1}^M \nabla f_m(\hat{\mathbf{w}}_m^{(t)}) \right\|^2 + \frac{\eta^2 L \sigma^2}{Mb^{(t)}} \\
 & = -\frac{\eta}{2M} \sum_{m=1}^M (\|\nabla F(\tilde{\mathbf{w}}^{(t)})\|^2 + L^2 \|\tilde{\mathbf{w}}^{(t)} - \hat{\mathbf{w}}_m^{(t)}\|^2) + \frac{2\eta^2 L - \eta}{2M} \sum_{m=1}^M \|\nabla f_m(\hat{\mathbf{w}}_m^{(t)})\|^2 + \frac{\eta^2 L \sigma^2}{Mb^{(t)}} + \frac{\eta L^2}{M} \|\tilde{\mathbf{w}}^{(t)} - \hat{\mathbf{w}}_m^{(t)}\|^2
 \end{aligned} \tag{2}$$

where (a) follows by applying two basic inequalities  $\langle \mathbf{a}, \mathbf{a} \rangle \leq 1/2 \|\mathbf{a}\|^2 + 1/2 \|\mathbf{b}\|^2$  and  $\mathbb{E}[\|X\|^2] = \mathbb{E}[\|X - \mathbb{E}[X]\|^2] + \|\mathbb{E}[X]\|^2$ ; (b) follows from the Lipschitz continuity of the gradient of local functions. The first term in (2) can be bounded in terms of  $\|\nabla f_m(\hat{\mathbf{w}}_m^{(t)})\|^2$  as follows:

$$\begin{aligned}
 \|\nabla f_m(\hat{\mathbf{w}}_m^{(t)})\|^2 & \leq 2\|\nabla f_m(\hat{\mathbf{w}}_m^{(t)}) - \nabla F(\tilde{\mathbf{w}}^{(t)})\|^2 + 2\|\nabla F(\tilde{\mathbf{w}}^{(t)})\|^2 \\
 & \leq 2L^2 \|\hat{\mathbf{w}}_m^{(t)} - \tilde{\mathbf{w}}^{(t)}\|^2 + \|\nabla F(\tilde{\mathbf{w}}^{(t)})\|^2
 \end{aligned} \tag{3}$$

Using  $\eta \leq \frac{1}{2L}$  and rearranging the terms in (2), we have

$$\frac{\eta}{4M} \sum_{m=1}^M \|\nabla f_m(\hat{\mathbf{w}}_m^{(t)})\|^2 \leq F(\tilde{\mathbf{w}}^{(t)}) - \mathbb{E}[F(\tilde{\mathbf{w}}^{(t+1)})] + \frac{\eta^2 L \sigma^2}{Mb^{(t)}} + \frac{\eta L^2}{M} \|\tilde{\mathbf{w}}^{(t)} - \hat{\mathbf{w}}_m^{(t)}\|^2 \tag{4}$$

Taking expectation with respect to the entire process and using the basic inequality  $\|\mathbf{a} + \mathbf{b}\|^2 \leq 2\|\mathbf{a}\|^2 + 2\|\mathbf{b}\|^2$  gives

$$\begin{aligned} & \frac{\eta}{4M} \sum_{m=1}^M \mathbb{E}[\|\nabla f_m(\hat{\mathbf{w}}_m^{(t)})\|^2] \\ & \leq \mathbb{E}[F(\tilde{\mathbf{w}}^{(t)})] - \mathbb{E}[F(\tilde{\mathbf{w}}^{(t+1)})] + \frac{\eta^2 L \sigma^2}{M b^{(t)}} + 2\eta L^2 \mathbb{E}[\|\tilde{\mathbf{w}}^{(t)} - \hat{\mathbf{w}}^{(t)}\|^2] + \frac{2\eta L^2}{M} \sum_{m=1}^M \mathbb{E}[\|\hat{\mathbf{w}}^{(t)} - \hat{\mathbf{w}}_m^{(t)}\|^2] \\ & \stackrel{(a)}{\leq} \mathbb{E}[F(\tilde{\mathbf{w}}^{(t)})] - \mathbb{E}[F(\tilde{\mathbf{w}}^{(t+1)})] + \frac{2\eta^2 L \sigma^2}{M \rho^t b^{(0)}} + 2\eta L^2 \mathbb{E}[\|\tilde{\mathbf{w}}^{(t)} - \hat{\mathbf{w}}^{(t)}\|^2] + \frac{2\eta L^2}{M} \sum_{m=1}^M \mathbb{E}[\|\hat{\mathbf{w}}^{(t)} - \hat{\mathbf{w}}_m^{(t)}\|^2], \end{aligned} \quad (5)$$

where (a) follows by recalling  $b^{(t)} = \lfloor \rho^t b^{(0)} \rfloor$  and noting  $\lfloor x \rfloor > x/2$  as long as  $x \geq 2$ .

Now we give three important lemmas where the first two are borrowed from [1] and the last one is proved in the following.

**Lemma 1 (Memory [1]):** The accumulated error captures the distance between the true sequence and virtual sequence. That is

$$\hat{\mathbf{w}}^{(t)} - \tilde{\mathbf{w}}^{(t)} = \frac{1}{M} \sum_{m=1}^M \mathbf{e}_m^{(t)} \quad (6)$$

**Lemma 2 (Contracting Deviation of Local Sequences [1]):** The deviation of the local sequences is bounded by

$$\frac{1}{M} \sum_{m=1}^M \mathbb{E}[\|\hat{\mathbf{w}}^{(t)} - \hat{\mathbf{w}}_m^{(t)}\|^2] \leq \eta^2 G^2 H^2 \quad (7)$$

**Lemma 3 (Bounded Memory):** For worker  $m$  who synchronizes with the server every  $H$  local iterations, we have

$$\mathbb{E}[\|\mathbf{e}_m^{(t)}\|^2] \leq 4\delta_m^2 \eta^2 G^2 H^2 \quad (8)$$

*Proof:* Note that Algorithm ?? average the gradients every  $H$  iterations between which the accumulated error  $\mathbf{e}_m^{(t)}$  at any participant  $m$  and the global parameter vector  $\mathbf{w}^{(t)}$  keep unchanged. For ease of presentation, we assume that  $T$  is an integer multiple of  $H$ . Let  $\mathcal{I}_T = \{t_1, t_2, \dots, t_{T/H} = T\}$  be the aggregation indices satisfying  $t_{i+1} - t_i = H$ . For every  $m \in \mathcal{M}$ , we have

$$\begin{aligned} \mathbb{E}[\|\mathbf{e}_m^{(t_{i+1})}\|^2] &= \mathbb{E}[\|\mathbf{e}_m^{(t_{i+1}-1)} + \mathbf{w}^{(t_{i+1}-1)} - \hat{\mathbf{w}}_m^{(t_{i+1}-\frac{1}{2})} - \mathbf{g}_m^{(t_{i+1}-1)}\|^2] \\ &\stackrel{(a)}{\leq} (1 - \frac{1}{\delta_m}) \mathbb{E}[\|\mathbf{e}_m^{(t_{i+1}-1)} + \mathbf{w}^{(t_{i+1}-1)} - \hat{\mathbf{w}}_m^{(t_{i+1}-\frac{1}{2})}\|^2] \\ &\stackrel{(b)}{=} (1 - \frac{1}{\delta_m}) \mathbb{E}[\|\mathbf{e}_m^{(t_i)} + \hat{\mathbf{w}}_m^{(t_i)} - \hat{\mathbf{w}}_m^{(t_{i+1}-\frac{1}{2})}\|^2]. \end{aligned} \quad (9)$$

Here (a) is due to the contraction property of  $\text{Top}_k(\mathbf{x})$  operator [2], that is  $\mathbb{E}[\|\mathbf{x} - \text{Top}_k(\mathbf{x})\|^2] \leq (1 - k/d)\|\mathbf{x}\|^2$ ,  $\forall \mathbf{x} \in \mathbb{R}^d$ . In (b), we use  $\mathbf{e}_m^{(t_{i+1}-1)} = \mathbf{e}_m^{(t_i)}$  and  $\mathbf{w}^{(t_{i+1}-1)} = \mathbf{w}^{(t_i)} = \hat{\mathbf{w}}_m^{(t_i)}$  that always hold. Since the inequality  $\|\mathbf{a} + \mathbf{b}\|^2 \leq (1 + \tau)\|\mathbf{a}\|^2 + (1 + \frac{1}{\tau})\|\mathbf{b}\|^2$  holds for every  $\tau \geq 0$ , we take any  $p > 1$  and transform (9) as follows

$$\begin{aligned} \mathbb{E}[\|\mathbf{e}_m^{(t_{i+1})}\|^2] &\leq (1 - \frac{1}{\delta_m}) \left\{ (1 + \frac{(p-1)}{p\delta_m}) \mathbb{E}[\|\mathbf{e}_m^{(t_i)}\|^2] + (1 + \frac{p\delta_m}{(p-1)}) \mathbb{E}[\|\hat{\mathbf{w}}_m^{(t_i)} - \hat{\mathbf{w}}_m^{(t_{i+1}-\frac{1}{2})}\|^2] \right\} \\ &\leq (1 - \frac{1}{p\delta_m}) \mathbb{E}[\|\mathbf{e}_m^{(t_i)}\|^2] + \frac{p(\delta_m^2 - 1)}{(p-1)\delta_m} \mathbb{E}[\|\sum_{j=t_i}^{t_{i+1}-1} \eta \nabla f_m(\mathbf{w}_m^{(j)}; \mathcal{D}_m^{(j)})\|^2] \\ &\stackrel{(a)}{\leq} (1 - \frac{1}{p\delta_m}) \mathbb{E}[\|\mathbf{e}_m^{(t_i)}\|^2] + \frac{p(\delta_m^2 - 1)}{(p-1)\delta_m} \eta^2 G^2 H^2, \end{aligned} \quad (10)$$

where (a) follows from Assumption 1. Iterating the above inequality from  $i = 0 \rightarrow l$  where  $l = T/H$  yields:

$$\mathbb{E}[\|\mathbf{e}_m^{(t_{i+1})}\|^2] \leq \frac{p(\delta_m^2 - 1)}{(p-1)\delta_m} \eta^2 G^2 H^2 \sum_{j=1}^l (1 - \frac{1}{p\delta_m})^{l-j} \stackrel{(a)}{\leq} \frac{p^2(\delta_m^2 - 1)}{p-1} \eta^2 G^2 H^2 \stackrel{(b)}{\leq} 4\delta_m^2 \eta^2 G^2 H^2, \quad (11)$$

where (a) is by the fact that  $\sum_{j=1}^l (1 - 1/p\delta_m)^{l-j} \leq \sum_{j \geq 0} (1 - 1/p\delta_m)^j = p\delta_m$ , and (b) is by plugging  $p = 2$ . Note the the right-hand-side does not depend on  $t$ , i.e., for every  $t = 0, 1, \dots, T$ , the following holds:

$$\mathbb{E}[\|\mathbf{e}_m^{(t)}\|^2] \leq 4\delta_m^2 \eta^2 G^2 H^2. \quad (12)$$

■

Lemma 1 and Lemma 3 together imply:

$$\mathbb{E}[\|\hat{\mathbf{w}}^{(t)} - \tilde{\mathbf{w}}^{(t)}\|^2] \leq \frac{4\eta^2 G^2 H^2}{M} \sum_{m=1}^M \delta_m^2. \quad (13)$$

Applying Lemma 2 and (13) into (5), we get

$$\frac{\eta}{4M} \sum_{m=1}^M \mathbb{E}[\|\nabla f_m(\hat{\mathbf{w}}_m^{(t)})\|^2] \leq \mathbb{E}[F(\tilde{\mathbf{w}}^{(t)})] - \mathbb{E}[F(\tilde{\mathbf{w}}^{(t+1)})] + \frac{2\eta^2 L \sigma^2}{M \rho^t b^{(0)}} + \frac{8\eta^3 L^2 G^2 H^2}{M} \sum_{m=1}^M \delta_m^2 + 2\eta^3 L^2 G^2 H^2, \quad (14)$$

Recursively applying the above inequality from  $t = 0$  to  $t = T - 1$  yields

$$\begin{aligned} \frac{1}{4MT} \sum_{t=0}^{T-1} \sum_{m=1}^M \mathbb{E}[\|\nabla f_m(\hat{\mathbf{w}}_m^{(t)})\|^2] &\leq \frac{\mathbb{E}[F(\tilde{\mathbf{w}}^{(0)})] - F^*}{\eta T} + \frac{2\eta L \sigma^2}{MT b^{(0)}} \sum_{t=0}^{T-1} \frac{1}{\rho^t} + \frac{8\eta^2 L^2 G^2 H^2}{M} \sum_{m=1}^M \delta_m^2 + 2\eta^2 L^2 G^2 H^2 \\ &\stackrel{(a)}{\leq} \frac{\mathbb{E}[F(\tilde{\mathbf{w}}^{(0)})] - F^*}{\eta T} + \frac{2\eta \rho L \sigma^2}{(\rho - 1)MT b^{(0)}} + \frac{8\eta^2 L^2 G^2 H^2}{M} \sum_{m=1}^M \delta_m^2 + 2\eta^2 L^2 G^2 H^2, \end{aligned} \quad (15)$$

where (a) follows by simplifying the partial sum of geometric series and noting that  $0 < \frac{1}{\rho} < 1$ . Let  $\mathbf{z}_T$  be a random variable sampled from  $\{\hat{\mathbf{w}}_m^{(t)}\}$  with probability  $\Pr[\mathbf{z}_T = \hat{\mathbf{w}}_m^{(t)}] = \frac{1}{MT}$ . By taking  $\delta = \sqrt{\frac{1}{M} \sum_{m=1}^M \delta_m^2}$  and  $\eta = \frac{\theta \sqrt{M}}{\sqrt{T}}$  (where  $\theta$  is a constant satisfying  $\frac{\theta \sqrt{M}}{\sqrt{T}} \leq \frac{1}{2L}$ ), we have

$$\mathbb{E}[\|\mathbf{z}_T\|^2] = \frac{1}{MT} \sum_{t=0}^{T-1} \sum_{m=1}^M \mathbb{E}[\|\nabla f_m(\hat{\mathbf{w}}_m^{(t)})\|^2] \leq \frac{4(\mathbb{E}[F(\mathbf{w}^{(0)})] - F^*)}{\theta \sqrt{MT}} + \frac{8\rho \theta L \sigma^2}{(\rho - 1)M b^{(0)} \sqrt{MT}^{3/2}} + (4\delta^2 + 1) \frac{8M\theta^2 L^2 G^2 H^2}{T}, \quad (16)$$

Until now we complete the proof of Theorem 1.

## REFERENCES

- [1] D. Basu, D. Data, C. Karakus, and S. Diggavi, "Qsparse-local-SGD: Distributed SGD with quantization, sparsification and local computations," in *Proc. of Advances in Neural Information Processing Systems (NIPS)*, Vancouver, Canada, December 2019.
- [2] S. U. Stich, J.-B. Cordonnier, and M. Jaggi, "Sparsified SGD with memory," in *Proc. of Advances in Neural Information Processing Systems (NIPS)*, Vancouver, Canada, December 2018.