

Checklist

1. For all authors...
 - (a) Do the main claims made in the abstract and introduction accurately reflect the paper's contributions and scope? [Yes]
 - (b) Did you describe the limitations of your work? [Yes]
 - (c) Did you discuss any potential negative societal impacts of your work? [N/A]
 - (d) Have you read the ethics review guidelines and ensured that your paper conforms to them? [Yes]
2. If you are including theoretical results...
 - (a) Did you state the full set of assumptions of all theoretical results? [Yes]
 - (b) Did you include complete proofs of all theoretical results? [Yes]
3. If you ran experiments...
 - (a) Did you include the code, data, and instructions needed to reproduce the main experimental results (either in the supplemental material or as a URL)? [No]
 - (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they were chosen)? [Yes]
 - (c) Did you report error bars (e.g., with respect to the random seed after running experiments multiple times)? [Yes]
 - (d) Did you include the total amount of compute and the type of resources used (e.g., type of GPUs, internal cluster, or cloud provider)? [No]
4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
 - (a) If your work uses existing assets, did you cite the creators? [Yes]
 - (b) Did you mention the license of the assets? [N/A]
 - (c) Did you include any new assets either in the supplemental material or as a URL? [N/A]
 - (d) Did you discuss whether and how consent was obtained from people whose data you're using/curating? [N/A]
 - (e) Did you discuss whether the data you are using/curating contains personally identifiable information or offensive content? [N/A]
5. If you used crowdsourcing or conducted research with human subjects...
 - (a) Did you include the full text of instructions given to participants and screenshots, if applicable? [N/A]
 - (b) Did you describe any potential participant risks, with links to Institutional Review Board (IRB) approvals, if applicable? [N/A]
 - (c) Did you include the estimated hourly wage paid to participants and the total amount spent on participant compensation? [N/A]

A Additional Related Literature and Comparison

In the first category, where the communication complexities are measured by *rounds* of communications (where in each round real-valued vectors get exchanged), recent works [Kovalev et al., 2021, 2022] provides lower bounds for a sum of smooth and strongly convex functions over time-varying networks, and for strongly monotone variational inequality problems in a stochastic (finite-sum) setting. Other related works include [Scaman et al., 2017, 2018, Arjevani and Shamir, 2015]. In addition, lower bounds for non-convex problems are considered in [Sun and Hong, 2019, Lu and De Sa, 2021]. Further, there are many works that derive rounds of communication upper bounds for decentralized and federated learning algorithms, see, e.g., [Stich and Karimireddy, 2019, Stich et al., 2018, Patel and Dieuleveut, 2019, Yu et al., 2019, Wang and Joshi, 2018, Gorbunov et al., 2021].

A recent work [Gorbunov et al., 2021] has analyzed the communication efficient algorithm to solve functions that satisfy PL conditions. However, it has been focused on analyzing the number of communication rounds needed to achieve certain ϵ optimal solution, while the current paper is focused on finding the minimum bits to be communicated. Therefore, although the two works are both about developing communication efficient algorithms for PL functions, the bounds obtained in these works represent different physical quantities, thus cannot be directly compared.

B Communication Complexity Lower Bounds

Let us introduce the so-called equality problem Vempala et al. [2020], denoted as EQUAL.

Definition 5 (Equality Problem). *Consider a set of K agents, each is given an input $c_k \in \{0, 1\}^m$; m is the length of the binary input. Then, the EQUAL problem is defined as follows:*

$$\text{EQUAL}_m(c_1, \dots, c_K) = \begin{cases} 1, & \text{if } c_1 = \dots = c_K \\ 0, & \text{otherwise} \end{cases}$$

For any deterministic algorithms, the communication complexity lower bound of $\text{EQUAL}_m(c_1, \dots, c_K)$ is $\Omega(Km)$ [Vempala et al., 2020, Thm 3.5].

Further, the intermediate steps required to derive the lower bounds involve packing arguments. Therefore, we provide below a lower bound for the maximum number of points that we can pack into a compact set $[0, 1]^n$, such that the distance between each pair of points is at least δ .

Definition 6 (Packing Problem). *We define the following:*

- For a given $\delta > 0$ we define the set $S(\delta) \subseteq [0, 1]^n$ such that $\|x - y\| > \delta, \forall x, y \in S(\delta)$.
- Assuming that $|S(\delta)| \geq 2^m$, we define a function $h : \{0, 1\}^m \rightarrow S(\delta)$. For $u, v \in \{0, 1\}^m$ it holds that $u \neq v \Leftrightarrow h(u) \neq h(v)$.

Lemma B.1. ([Korhonen and Alistarh, 2021, Lemma 2]) *For a set $S(\delta) \subseteq [0, 1]^n$ defined in Def. 6, it holds that $|S(\delta)| \geq \left(\frac{\sqrt{2n}}{\sqrt{\pi\epsilon\delta}}\right)^n$.*

Next, we repeat here for completeness the Assumptions that the local functions in the Distributed PL problem class satisfy (Def. 1).

Assumption 6. *The local objective functions satisfy:*

$$2\mu_k \cdot (f_k(\theta) - f_k(\theta_{(k)}^*)) \leq \|\nabla f_k(\theta)\|^2, \quad \forall \theta, \quad \forall k,$$

where $\theta_{(k)}^*$ is a global minimum of $f_k(\cdot)$; μ_k 's some positive constants.

Assumption 7. *There exists positive constants L_k 's and L such that:*

$$\|\nabla f_k(\theta) - \nabla f_k(\theta')\| \leq L_k \|\theta - \theta'\|, \quad \|\nabla f(\theta) - \nabla f(\theta')\| \leq L \|\theta - \theta'\|, \quad \forall \theta, \theta', \quad \forall k.$$

In the Lemma below we show that the function instance we are going to use in Theorem 3.1 satisfies the PL condition and both Assumptions 6 and 7.

Lemma B.2. *The function $f_k(\theta) = \frac{1}{2}\|\theta - c\|_2^2 + \sin^2(\theta_k - c_k), c \in \mathbb{R}^D$ (where with c_k, θ_k we denote the k th component of vectors c, θ , respectively) satisfies the PL condition with $\mu = \frac{1}{8}$. That is, it holds that $f_k(\theta) - f_k(\theta^*) \leq \frac{1}{2\mu}\|\nabla f_k(\theta)\|^2, \forall \theta \in \mathbb{R}^D$, where θ^* is the global minimum of $f_k(\cdot)$.*

Remark 1. The constructed function $f_k(\theta)$ and the respective sum across nodes $\sum_{k=1}^K f_k(\theta)$ satisfy Assumptions 6, 7. Specifically, Assumption 6 follows trivially from Lemma B.2. In addition, it can be shown that $f_k(\theta)$ and $\sum_{k=1}^K f_k(\theta)$ have bounded Hessians, which implies the Lipschitz gradient property of Assumption 7.

Proof of Lemma B.2. To begin with we are going to show that the (one dimensional) function

$$\tilde{f}_k(\theta_k) = \frac{1}{2} (\theta_k - c_k)^2 + \sin^2 (\theta_k - c_k),$$

where $\theta_k, c_k \in \mathbb{R}$, satisfies the PL condition with $\mu = \frac{1}{8}$; notice that $\tilde{f}_k(\theta_k)$ has a unique minimum with value 0 attained at $\theta_k = c_k$, and we have that $\tilde{f}'_k(\theta_k) = \theta_k - c_k + \sin(2(\theta_k - c_k))$. Then, we are going to use this result to prove the PL condition for the function $f_k(\theta)$.

So, let us define the function

$$g(\theta_k) = 4 [\theta_k - c_k + \sin(2(\theta_k - c_k))]^2 - \frac{1}{2} (\theta_k - c_k)^2 - \sin^2 (\theta_k - c_k).$$

The gradient of the above function is given by

$$\begin{aligned} g'(\theta_k) &= 8 [\theta_k - c_k + \sin(2(\theta_k - c_k))] [1 + 2 \cos(2(\theta_k - c_k))] - (\theta_k - c_k) - \sin(2(\theta_k - c_k)) \\ &= [\theta_k - c_k + \sin(2(\theta_k - c_k))] [7 + 16 \cos(2(\theta_k - c_k))]. \end{aligned}$$

In order to show that $\tilde{f}_k(\theta_k)$ satisfies the PL property it suffices to prove that $g(\theta_k) \geq 0, \forall \theta_k \in \mathbb{R}$. Then, it will hold that

$$4 [\theta_k - c_k + \sin(2(\theta_k - c_k))]^2 \geq \frac{1}{2} (\theta_k - c_k)^2 - \sin^2 (\theta_k - c_k), \forall \theta_k \in \mathbb{R},$$

that is the PL condition of $\tilde{f}_k(\theta_k)$ will be satisfied.

First, notice that

$$\begin{aligned} g(\theta_k) &= 4 (\theta_k - c_k)^2 + 8 (\theta_k - c_k) \sin(2(\theta_k - c_k)) + 4 \sin^2(2(\theta_k - c_k)) - \frac{1}{2} (\theta_k - c_k)^2 - \sin^2 (\theta_k - c_k) \\ &\geq \frac{7}{2} (\theta_k - c_k)^2 + 8 (\theta_k - c_k) \sin(2(\theta_k - c_k)) - 1. \end{aligned}$$

It is clear from the above expression that the term $\frac{7}{2} (\theta_k - c_k)^2$ dominates the value of the objective for large values of $|\theta_k - c_k|$. Therefore, the objective does not become unbounded below.

Secondly, consider the stationary points of $g(\theta_k)$. Since the objective does not become unbounded below the only possible global minima of $g(\theta_k)$ are its stationary points. We are going to show that the values of the objective at those points is non-negative, effectively proving that $g(\theta_k) \geq 0, \forall \theta_k \in \mathbb{R}$. The stationary points of g are defined by the following expressions:

- $\sin(2(\theta_k - c_k)) = -(\theta_k - c_k)$

Notice that in the interval $\theta_k \in (c_k, \frac{\pi}{2} + c_k]$ it holds that $\sin[2(\theta_k - c_k)] \geq 0$ but $-(\theta_k - c_k) < 0$. Similarly, for $\theta_k \in [-\frac{\pi}{2} + c_k, c_k)$ it holds that $\sin[2(\theta_k - c_k)] \leq 0$ but $-(\theta_k - c_k) > 0$. Also, for $\theta_k \notin [-\frac{\pi}{2} + c_k, c_k) \cup (c_k, \frac{\pi}{2} + c_k] \cup \{c_k\}$ it holds that $|\theta_k - c_k| \geq \frac{\pi}{2} > 1$. Therefore, in all the above cases the equation of the stationary points does not have a solution (and thus there are no stationary points). Finally, note that the equation has a trivial solution at $\theta_k = c_k$, which corresponds to the only stationary point we can get from this equation. For that point it holds that $g(c_k) = 0$.

- $\cos(2(\theta_k - c_k)) = -\frac{7}{16}$

From the above equations, and by using the proper trigonometric identities it follows that $\sin^2(2(\theta_k - c_k)) = \frac{207}{256}$, $\sin^2(\theta_k - c_k) = \frac{23}{32}$ and $\theta_k - c_k = \kappa\pi \pm \frac{1}{2} \arccos(-\frac{7}{16}) \approx \kappa\pi \pm \frac{2.02}{2}$. Then if we plug the above values into $g(\theta_k)$ we get

$$g(\theta_k) \approx \frac{7}{2} (\kappa\pi \pm 1.01)^2 + 8 (\kappa\pi \pm 1.01) \sin[2\kappa\pi \pm 2.02] + 4 \frac{207}{256} - \frac{23}{32}$$

$$\begin{aligned}
&= \frac{7}{2} (\kappa\pi \pm 1.01)^2 + 8 (\kappa\pi \pm 1.01) \sin [\pm 2.02] + \frac{644}{256} \\
&= (\kappa\pi \pm 1.01) \left[\frac{7}{2} (\kappa\pi \pm 1.01) + 8 \cdot (\pm 0.9) \right] + \frac{161}{64}.
\end{aligned}$$

It can be easily verified that $(\kappa\pi \pm 1.01) \left[\frac{7}{2} (\kappa\pi \pm 1.01) + 8 \cdot (\pm 0.9) \right] + \frac{161}{64} > 0$ for all $\kappa \in \mathbb{Z}$.

In conclusion for all the stationary points $\hat{\theta}_k$ it holds that $g(\hat{\theta}_k) \geq 0$. As a result, we can claim that $g(\theta_k) \geq 0, \forall \theta_k \in \mathbb{R}$, and the function $\tilde{f}(\theta_k)$ satisfies the PL property.

Next, consider the PL property (with $\mu = \frac{1}{8}$) for the function $f_k(\theta)$, that is

$$\frac{1}{2} \|\theta - c\|_2^2 + \sin^2(\theta_k - c_k) \leq 4 \|\theta - c + \sin(2(\theta_k - c_k)) e_k\|^2, \quad (9)$$

where e_k is a vector of all zeros except at index k .

We have shown above that

$$\frac{1}{2} (\theta_k - c_k)^2 + \sin^2(\theta_k - c_k) \leq 4 [\theta_k - c_k + \sin(2(\theta_k - c_k))]^2. \quad (10)$$

Also, it trivially holds that

$$\frac{1}{2} \sum_{i=1, i \neq k}^D (\theta_i - c_i)^2 \leq 4 \sum_{i=1, i \neq k}^D (\theta_i - c_i)^2. \quad (11)$$

Adding inequalities (10) and (11) we obtain condition (9), which ensures the PL property for the objective $f_k(\theta)$, with $\mu = \frac{1}{8}$. This completes the proof. \square

Proof of Theorem 3.1. To begin with, let $\pi \in \Pi(\epsilon)$ be an arbitrary protocol that solves the problems in \mathcal{C}_{pl} . Then, for any input (i.e., for any specific problem within the class \mathcal{C}_{pl}) the protocol returns an ϵ -approximate minimum $\tilde{\theta}$ as given in Def. 3.

Moreover, consider the set $S(\delta) \subseteq [0, 1]^D$ introduced in Def. 6 (where in place of n we have D) with $\delta = 2\sqrt{2}\epsilon$. Then, Lemma B.1 implies that

$$|S(\delta)| \geq \left(\frac{\sqrt{2D}}{2\sqrt{2\pi e\epsilon}} \right)^D = \left(\frac{1}{2\sqrt{\pi e}} \right)^D \left(\frac{D}{\epsilon} \right)^{D/2}.$$

Then, we set $m = \Theta(\log_2 |S(\delta)|) = \Theta(D \log(\frac{D}{\epsilon}))$, where the exact value of m is selected such that $|S(\delta)| > 2^m$ holds. As a result, under the assumption that $\frac{D}{\epsilon} = \Omega(1)$, it holds that $|S(\delta)| \geq 2$, and $m \geq 1$.

Next, we will show that protocol π also solves the EQUAL problem. That is, we will reduce every instance of EQUAL $_m(u_1, \dots, u_K)$ to a problem in \mathcal{C}_{pl} . Towards this end, for an arbitrary input $(u_1, \dots, u_K) \in \{0, 1\}^{K^m}$ we select the following function from \mathcal{C}_{pl} ,

$$f(\theta) = \sum_{k=1}^K f_k(\theta) \text{ with } f_k(\theta) = \frac{1}{2} \|\theta - h(u_k)\|_2^2 + \sin^2(\theta_k - (h(u_k))_k), \quad (12)$$

where $h(\cdot)$ is introduced in Def. 6, and $h(u_k) \in S(\delta) \subseteq [0, 1]^D$; $(h(u_k))_k, \theta_k$ denote the k th component of $h(u_k)$ and θ , respectively. It is shown in Lemma B.2 that the functions f_k 's that correspond to each node satisfy the PL condition, with $\mu_k = \frac{1}{8}$, for all k . We can also easily verify that the rest of the conditions in Assumption 6–7 are satisfied.

Then, we have the following cases:

- **Case 1** (equal inputs): It holds that $u_1 = \dots = u_K := u$. Thus, we have that $h(u_1) = \dots = h(u_K) = h(u)$ and the minimum of (12) is 0, attained at some point θ^* such that $\theta^* = h(u)$.

Also, protocol π returns an approximate minimum $f(\tilde{\theta})$ of (12), for which it holds that

$$f(\tilde{\theta}) \stackrel{(a)}{\leq} f(\theta^*) + \epsilon = \frac{1}{2} \sum_{k=1}^K \|\theta^* - h(u)\|_2^2 + \sin^2(\theta_k^* - (h(u))_k) + \epsilon \stackrel{(b)}{=} \epsilon,$$

where in (a) we used (3); and in (b) we exploited the fact that the minimum is 0.

- **Case 2** (inputs not equal): There exists a pair of nodes (i, j) such that $u_i \neq u_j$. As a result, $h(u_i) \neq h(u_j)$. Then, protocol π returns an approximate minimum $f(\tilde{\theta})$ of (12) for which it holds that

$$\begin{aligned} f(\tilde{\theta}) &\stackrel{(a)}{\geq} f(\theta^*) \\ &= f_i(\theta^*) + f_j(\theta^*) + \sum_{k=1, k \neq i, j}^K f_k(\theta^*) \\ &\geq \frac{1}{2} \|\theta^* - h_i(u_i)\|_2^2 + \frac{1}{2} \|\theta^* - h_j(u_j)\|_2^2 \\ &\stackrel{(b)}{>} \frac{1}{2} \left(\frac{\delta}{2}\right)^2 = \frac{\delta^2}{8} \stackrel{(c)}{=} \epsilon, \end{aligned}$$

where in (a) we used (3), in (b) we used the characteristic property (i.e., the minimum distance between two points is $\|x - y\| > \delta, \forall x, y \in S(\delta)$) of set $|S(\delta)|$ (from Def. 6), and in (c) we use the quantity $\delta = 2\sqrt{2}\epsilon$.

From the above analysis we see that if $f(\tilde{\theta}) \leq \epsilon$ then $\text{EQUAL}_m(u_1, \dots, u_K) = 1$. Otherwise, if we assume that $\text{EQUAL}_m(u_1, \dots, u_K) = 0$, then the analysis of case 2 implies that $f(\tilde{\theta}) > \epsilon$, a contradiction. Similarly, we can claim that if $f(\tilde{\theta}) > \epsilon$, then $\text{EQUAL}_m(u_1, \dots, u_K) = 0$. In the opposite case (i.e., if $f(\tilde{\theta}) > \epsilon \implies \text{EQUAL}_m(u_1, \dots, u_K) = 1$) we see from case 1 that $f(\tilde{\theta}) \leq \epsilon$, that is we reach a contradiction. In summary, we have that

$$\text{EQUAL}_m(u_1, \dots, u_K) = \begin{cases} 1, & \text{if } f(\tilde{\theta}) \leq \epsilon \\ 0, & \text{if } f(\tilde{\theta}) > \epsilon. \end{cases}$$

Finally, the fact that the communication complexity of $\text{EQUAL}_m(u_1, \dots, u_K)$ is $\Omega(Km)$, and the above reduction imply that $\Omega(Km) = \Omega(KD \log(\frac{D}{\epsilon}))$ is a lower bound for the communication complexity of \mathcal{C}_{pl} . \square

Proof of Theorem 3.2. To begin with, let $\pi \in \Pi(\epsilon)$ be an arbitrary protocol that solves the problems in \mathcal{C}_{op} . Then, for any input (i.e., for any specific problem within the class \mathcal{C}_{op}) the protocol returns an ϵ -approximate minimum $\tilde{\theta}$, as described in Def. 3.

Moreover, consider the set $S(\delta) \subseteq [0, 1]^N$ introduced in Def. 6 (where in place of n we have N) with $\delta = 2\sqrt{2}\epsilon$. Then, Lemma B.1 implies that

$$|S(\delta)| \geq \left(\frac{\sqrt{2N}}{2\sqrt{2\pi e\epsilon}}\right)^N = \left(\frac{1}{2\sqrt{\pi e}}\right)^N \left(\frac{N}{\epsilon}\right)^{N/2}.$$

Then, we set $m = \Theta(\log_2 |S(\delta)|) = \Theta(N \log(\frac{N}{\epsilon}))$, where the exact value of m is selected such that $|S(\delta)| > 2^m$ holds. As a result, under the assumption that $\frac{N}{\epsilon} = \Omega(1)$, it holds that $|S(\delta)| \geq 2$, and $m \geq 1$.

Now, we are going to show that protocol π also solves the EQUAL problem. That is, we are going to reduce every instance of $\text{EQUAL}_m(u_1, \dots, u_K)$ to a problem in \mathcal{C}_{op} . To be more precise, for an

arbitrary input $(u_1, \dots, u_K) \in \{0, 1\}^{Km}$ we select the following instance from \mathcal{C}_{op} ,

$$f(\theta) = \sum_{k=1}^K f_k(\theta) \text{ with } f_k(\theta) = \frac{1}{2} \|G(\theta) - h(u_k)\|_2^2, \quad (13)$$

where $h(\cdot)$ is introduced in Def. 6; and with $h(u_k) \in S \subseteq [0, 1]^N$.

Then, we have the following cases:

- **Case 1** (equal inputs): It holds that $u_1 = \dots = u_K := u$. Thus, we have that $h(u_1) = \dots = h(u_K) = h(u)$ and the minimum of (13) is 0, attained at some point θ^* such that $G(\theta^*) = h(u)$.

Also, protocol π returns an approximate minimum $f(\tilde{\theta})$ of (13) for which it holds that

$$f(\tilde{\theta}) \stackrel{(a)}{\leq} f(\theta^*) + \epsilon = \frac{1}{2} \sum_{k=1}^K \|G(\theta^*) - h(u)\|_2^2 + \epsilon \stackrel{(b)}{=} \epsilon,$$

where in (a) we used (3), and (b) follows from the fact that the minimum of (13) is 0 in this case (i.e., equal inputs).

- **Case 2** (inputs not equal): There exists a pair of nodes (i, j) such that $u_i \neq u_j$. As a result, $h(u_i) \neq h(u_j)$. Then, protocol π returns an approximate minimum $f(\tilde{\theta})$ of (13) for which it holds that:

$$\begin{aligned} f(\tilde{\theta}) &\stackrel{(a)}{\geq} f(\theta^*) \\ &= \frac{1}{2} \|G(\theta^*) - h(u_i)\|_2^2 + \frac{1}{2} \|G(\theta^*) - h(u_j)\|_2^2 + \frac{1}{2} \sum_{k=1, k \neq i, j}^K \|G(\theta^*) - h(u_k)\|_2^2 \\ &\geq \frac{1}{2} \|G(\theta^*) - h(u_i)\|_2^2 + \frac{1}{2} \|G(\theta^*) - h(u_j)\|_2^2 \\ &\stackrel{(b)}{>} \frac{1}{2} \left(\frac{\delta}{2}\right)^2 = \frac{\delta^2}{8} \stackrel{(c)}{=} \epsilon, \end{aligned}$$

where in (a) expression (3) is used, in (b) we used the characteristic property (i.e., the minimum distance between two points is $\|x - y\| > \delta, \forall x, y \in S(\delta)$) of set $|S(\delta)|$ (from Def. 6), and in (c) we use the quantity $\delta = 2\sqrt{2}\epsilon$.

From the above analysis we see that if $f(\tilde{\theta}) \leq \epsilon$ then $\text{EQUAL}_m(u_1, \dots, u_K) = 1$. Otherwise, if $\text{EQUAL}_m(u_1, \dots, u_K) = 0$, then the analysis of case 2 implies that $f(\tilde{\theta}) > \epsilon$, a contradiction. Similarly, we can claim that if $f(\tilde{\theta}) > \epsilon$, then $\text{EQUAL}_m(u_1, \dots, u_K) = 0$. In the opposite case (i.e., if $f(\tilde{\theta}) > \epsilon \implies \text{EQUAL}_m(u_1, \dots, u_K) = 1$) we see from case 1 that $f(\tilde{\theta}) \leq \epsilon$, that is we reach a contradiction. In summary, we have that

$$\text{EQUAL}_m(u_1, \dots, u_K) = \begin{cases} 1, & \text{if } f(\tilde{\theta}) \leq \epsilon \\ 0, & \text{if } f(\tilde{\theta}) > \epsilon. \end{cases}$$

Finally, the fact that the communication complexity of $\text{EQUAL}_m(u_1, \dots, u_K)$ is $\Omega(Km)$, and the above reduction imply that $\Omega(Km) = \Omega(KN \log(\frac{N}{\epsilon}))$ is a lower bound for the communication complexity of \mathcal{C}_{op} . \square

C Proof for Theorem 4.1

C.1 Proof of Lemma 4.2

Let us consider the function

$$\phi(\theta) = f(\theta) - \langle \nabla f(\theta^*), \theta \rangle.$$

Since $\nabla f(\theta^*) = 0$, it is easy to see $\phi(\theta) \geq 0$. It can be derived directly that $\phi(\theta^*) = 0$, which means $\phi(\cdot)$ can achieve the minimum at θ^* . So we know $\theta^* \in \text{argmin } \phi(\theta)$. Then we have the following inequality holds:

$$\begin{aligned}\phi(\theta^*) &\leq \phi\left(\theta - \frac{1}{L}\nabla\phi(\theta)\right) \\ &\stackrel{(i)}{=} f\left(\theta - \frac{1}{L}\nabla\phi(\theta)\right) \stackrel{(ii)}{=} f\left(\theta - \frac{1}{L}\nabla f(\theta)\right) \\ &\stackrel{(iii)}{\leq} f(\theta) - \frac{1}{2L}\|\nabla f(\theta)\|^2,\end{aligned}$$

where (i) performed one step of gradient decent; (ii) uses the fact that $\nabla f(\theta^*) = 0$; (iii) is because the Lipschitz gradient assumption. Then it follows directly

$$\|\nabla f(\theta)\|^2 \leq 2L \cdot f(\theta).$$

C.2 Proof for Theorem 4.1

First, let us state the sketch of the proof. Denote $g^t := \sum_{k=1}^K g_k^t$, $q^t := \sum_{k=1}^K q_k^t$.

Step1: We show that the loss function decreases linearly if all the agents update the averaged gradient in one step. That is, the following holds true.

$$f(\theta^t - \eta g^t) \leq (1 - \eta\mu)f(\theta^t).$$

Step2: We show by induction that for $t = 1, \dots$, the following inequalities hold true:

$$(1) f(\theta^t) \leq (\alpha)^t f(\theta^0), \text{ where } 0 < \alpha < 1, \quad (14)$$

$$(2) \|g^t - q^t\|_\infty \leq \frac{\gamma^{t-1}}{2^b - 1} = \tau\gamma^{t-1}, \text{ for some } \gamma^{t-1} > 0, \quad (15)$$

where $\{\gamma^t\}_{t=1}^\infty$ is a sequence of positive numbers.

The proof of **Step 1** is straightforward:

$$\begin{aligned}f(\theta^t - \eta g^t) &= f\left(\theta^t - \eta \sum_{k=1}^K g_k^t\right) \\ &\stackrel{(i)}{\leq} f(\theta^t) - \langle \nabla f(\theta^t), \eta \sum_{k=1}^K g_k^t \rangle + \frac{\eta^2 L}{2} \left\| \sum_{k=1}^K g_k^t \right\|^2 \\ &\stackrel{(ii)}{=} f(\theta^t) - \eta \|g^t\|^2 + \frac{\eta^2 L}{2} \|g^t\|^2 \\ &\stackrel{(iii)}{=} f(\theta^t) - \eta \|g^t\|^2 + \frac{\eta}{4L^2} \|g^t\|^2 \quad \left(\eta = \frac{1}{2L^{\frac{3}{2}}}\right) \\ &\stackrel{(iv)}{\leq} f(\theta^t) - \frac{1}{2}\eta \|g^t\|^2 \\ &\stackrel{(v)}{\leq} f(\theta^t) - \eta\mu f(\theta^t) \\ &= (1 - \eta\mu)f(\theta^t),\end{aligned} \quad (16)$$

where (i) is by Decent Lemma; (ii) uses the definition of g^t ; (iii) uses the choice of η ; (iv) is because $L > 1$; (v) uses Assumption 1.

Now we prove **Step 2**. To begin with, let us set $\eta = \frac{1}{2L^{\frac{3}{2}}}$, and set:

$$\alpha = 1 - \frac{\mu}{8L^{\frac{3}{2}}}, \quad \gamma^t = \sqrt{(\alpha)^{t+1} f(\theta^0)}, \quad \tau = \frac{1}{\sqrt{CD}}, \quad (17)$$

$$C = \max\left(\sqrt{\frac{16L}{\mu^2} + \frac{L^{\frac{3}{2}}}{\mu}}, 100\right), \quad b = \max\left(\log\left(\frac{1}{\tau} + 1\right), b_0\right).$$

First, let us verify that (14) holds true when $t = 1$. We have the following series of inequalities:

$$\begin{aligned}
f(\theta^1) &= f(\theta^0 - \eta q^0) - f(\theta^0 - \eta g^0) + f(\theta^0 - \eta g^0) \\
&\stackrel{(i)}{\leq} \eta \langle \nabla f(\theta^0 - \eta g^0), q^0 - g^0 \rangle + \frac{\eta^2 L}{2} \|q^0 - g^0\|^2 + f(\theta^0 - \eta g^0) \\
&\stackrel{(ii)}{\leq} \frac{\eta}{2\beta} \|\nabla f(\theta^0 - \eta g^0)\|^2 + 2\eta\beta \|q^0 - g^0\|^2 + \frac{\eta^2 L}{2} \|q^0 - g^0\|^2 + f(\theta^0 - \eta g^0) \\
&\stackrel{(iii)}{\leq} \frac{\eta L}{\beta} f(\theta^0 - \eta g^0) + (2\eta\beta + \frac{\eta^2 L}{2}) \|q^0 - g^0\|^2 + f(\theta^0 - \eta g^0) \\
&\stackrel{(iv)}{=} (\frac{\eta L}{\beta} + 1)(1 - \eta\mu) f(\theta^0) + (2\eta\beta + \frac{\eta^2 L}{2}) \|q^0 - g^0\|^2 \\
&\stackrel{(v)}{\leq} (\frac{\eta L}{\beta} + 1)(1 - \eta\mu) f(\theta^0) + (2\eta\beta + \frac{\eta^2 L}{2}) D^2 \|q^0 - g^0\|_\infty^2 \\
&\stackrel{(vi)}{\leq} (\frac{\eta L}{\beta} + 1)(1 - \eta\mu) f(\theta^0) + (2\eta\beta + \frac{\eta^2 L}{2}) D^2 (\sum_{k=1}^K \|q_k^0 - g_k^0\|_\infty)^2 \\
&\stackrel{(vii)}{\leq} (\frac{\eta L}{\beta} + 1)(1 - \eta\mu) f(\theta^0) + (2\eta\beta + \frac{\eta^2 L}{2}) D^2 \sum_{k=1}^K K \|q_k^0 - g_k^0\|_\infty^2 \\
&\stackrel{(viii)}{\leq} (\frac{\eta L}{\beta} + 1)(1 - \eta\mu) f(\theta^0) + (2\eta\beta + \frac{\eta^2 L}{2}) \frac{f(\theta^0)}{C^2 D^4} \\
&\stackrel{(ix)}{=} (1 - \frac{1}{2}\eta\mu) f(\theta^0) + (\frac{4\eta(1 - \eta\mu)L}{\mu} + \frac{\eta^2 L}{2}) \frac{f(\theta^0)}{C^2 D^4} \quad \beta = \frac{2(1 - \eta\mu)L}{\mu} \\
&= (1 - \frac{1}{2}\eta\mu + (\frac{4\eta(1 - \eta\mu)L}{\mu} + \frac{\eta^2 L}{2}) / C^2 D^4) f(\theta^0)
\end{aligned}$$

where (i) comes from the Decent Lemma; (ii) uses the Young's inequality with constant β ; (iii) uses Assumption 1; (iv) is from (16); (v) uses the relationship between ℓ_2 and ℓ_∞ norm; (vi) uses the triangle inequality; (vii) uses the Cauchy-Schwartz inequality; (viii) uses the initialization condition $\|q_k^0 - g_k^0\|_\infty \leq \frac{\sqrt{f_k(\theta^0)}}{CD^3\sqrt{K}}$; (ix) uses the choice of β .

Now let us define the quantization at initialization. For each entry in q_k^0 , denoted as $(g_k^0)_j$, we consider the interval $[\lfloor (g_k^0)_j \rfloor, \lceil (g_k^0)_j \rceil]$, which is constructed by the closest integers. We use b_0 bits to make the grid, and quantize each element of the vector by the closest point on the grid. It is clear that the quantization error for each element is at most $\frac{1}{2^{b_0}-1}$. So we set

$$\frac{1}{2^{b_0}-1} = \frac{\sqrt{f(\theta^0)}}{CD^3\sqrt{K}},$$

or equivalently, $b_0 = \log(\frac{CD^3\sqrt{K}}{\sqrt{f(\theta^0)}} + 1)$.

Now since we have chosen $C \geq \frac{16L}{\mu^2} + \frac{L^{\frac{3}{2}}}{\mu}$, we can bound the coefficient in front of $f(\theta^0)$ as:

$$\begin{aligned}
&1 - \frac{1}{2}\eta\mu + (\frac{4\eta(1 - \eta\mu)L}{\mu} + \frac{\eta^2 L}{2}) / C^2 D^4 \\
&\stackrel{(i)}{\leq} 1 - \frac{1}{2}\eta\mu + (\frac{2(1 - \eta\mu)}{\sqrt{L}\mu} + \frac{1}{8L^2}) / C^2 \\
&\stackrel{(ii)}{\leq} 1 - \frac{1}{2}\eta\mu + (\frac{2}{\sqrt{L}\mu} + \frac{1}{8}) / C^2 \\
&\stackrel{(iii)}{\leq} 1 - \frac{1}{4}\eta\mu \tag{18}
\end{aligned}$$

where in (i) we plugged in the choice of η and $D \geq 1$; (ii) is because we have assumed that $L > 1$ and $D \geq 1$; (iii) uses the choice of C in (17). It follows $f(\theta^1) \leq (1 - \frac{\mu}{8L^{\frac{3}{2}}}) f(\theta^0)$. Thus, (14) holds for $t = 1$.

Second, let us analyze (15) for $t = 1$. Observe that:

$$\begin{aligned}
\|g^1 - q^0\|_\infty &\leq \|g^1 - g^0\| + \|g^0 - q^0\|_\infty \stackrel{(i)}{\leq} \eta L \|q^0\| + \|g^0 - q^0\|_\infty \\
&\stackrel{(ii)}{\leq} \eta L \|g^0\| + \eta L \|q^0 - g^0\| + \|q^0 - g^0\|_\infty \\
&\stackrel{(iii)}{\leq} \eta L \|g^0\| + (1 + \eta LD) \|q^0 - g^0\|_\infty \\
&\stackrel{(iv)}{\leq} \eta L \|g^0\| + (1 + \eta LD) \sum_{k=1}^K \frac{\sqrt{f_k(\theta^0)}}{CD^3 \sqrt{K}} \\
&\stackrel{(v)}{\leq} \eta L \sqrt{2L \cdot f(\theta^0)} + (1 + \eta LD) \frac{\sqrt{f(\theta^0)}}{CD^3} \\
&\stackrel{(vi)}{\leq} \frac{\sqrt{2}}{2} \sqrt{f(\theta^0)} + 2 \cdot \frac{\sqrt{f(\theta^0)}}{100} \stackrel{(vii)}{\leq} \sqrt{\alpha f(\theta^0)},
\end{aligned}$$

where (i) uses Assumption 2; (ii) is from triangle inequality; (iii) uses the relationship between ℓ_2 and ℓ_∞ norm; (iv) uses the condition $\|q_k^0 - g_k^0\|_\infty \leq \frac{\sqrt{f_k(\theta^0)}}{CD^3 \sqrt{K}}$; (v) uses Assumption 1 and the Cauchy-Schwartz inequality; (vi) plugs in the choice of stepsize and the fact that $C \geq 100$; (v) compares the left side with definition of α in (17). By Lemma 4.1, it follows that:

$$\|q^1 - g^1\|_\infty = \|\text{quant}(g^1, q^0, r, b) - g^1\|_\infty \leq \frac{\gamma^0}{2^b - 1}.$$

Next, let us use induction to prove (14) and (15). First, we analyze the decent of $f(\theta^t)$:

$$\begin{aligned}
f(\theta^{t+1}) &= f(\theta^t - \eta q^t) - f(\theta^t - \eta g^t) + f(\theta^t - \eta g^t) \\
&\stackrel{(i)}{\leq} \eta \langle \nabla f(\theta^t - \eta g^t), q^t - g^t \rangle + \frac{\eta^2 L}{2} \|q^t - g^t\|^2 + f(\theta^t - \eta g^t) \\
&\stackrel{(ii)}{\leq} \frac{\eta}{2\beta} \|\nabla f(\theta^t - \eta g^t)\|^2 + 2\eta\beta \|q^t - g^t\|^2 + \frac{\eta^2 L}{2} \|q^t - g^t\|^2 + f(\theta^t - \eta g^t) \\
&\stackrel{(iii)}{\leq} \frac{\eta L}{\beta} f(\theta^t - \eta g^t) + (2\eta\beta + \frac{\eta^2 L}{2}) \|q^t - g^t\|^2 + f(\theta^t - \eta g^t) \\
&\stackrel{(iv)}{=} (\frac{\eta L}{\beta} + 1)(1 - \eta\mu) f(\theta^t) + (2\eta\beta + \frac{\eta^2 L}{2}) \tau^2 (\gamma^{t-1})^2 \\
&\stackrel{(v)}{\leq} (\frac{\eta L}{\beta} + 1)(1 - \eta\mu) (\alpha)^t f(\theta^0) + (2\eta\beta + \frac{\eta^2 L}{2}) \tau^2 (\gamma^{t-1})^2 \\
&\stackrel{(vi)}{\leq} (\frac{\eta L}{\beta} + 1)(1 - \eta\mu) (\alpha)^t f(\theta^0) + (2\eta\beta + \frac{\eta^2 L}{2}) \tau^2 (\alpha)^t f(\theta^0) \\
&= \left(1 - \frac{1}{2}\eta\mu + \left(\frac{4\eta(1 - \eta\mu)L}{\mu} + \frac{\eta^2 L}{2}\right) \tau^2\right) (\alpha)^t f(\theta^0), \tag{19}
\end{aligned}$$

where (i) comes from the Decent Lemma; (ii) uses the Young's inequality with constant β ; (iii) uses Assumption 1; (iv) is from (16) and induction assumption (15); (v) uses the induction assumption (14); (vi) uses the definition of γ^{t-1} . Finally, set $\beta = \frac{2(1 - \eta\mu)L}{\mu}$, we will get the last equality. Next, let us analyze the coefficient in the expression (19):

$$\begin{aligned}
1 - \frac{1}{2}\eta\mu + \left(\frac{4\eta(1 - \eta\mu)L}{\mu} + \frac{\eta^2 L}{2}\right) \tau^2 &\stackrel{(i)}{=} 1 - \frac{1}{2}\eta\mu + \left(\frac{2(1 - \eta\mu)}{\sqrt{L}\mu} + \frac{1}{8L^2}\right) \tau^2 \\
&\stackrel{(ii)}{\leq} 1 - \frac{1}{2}\eta\mu + \left(\frac{2}{\sqrt{L}\mu} + \frac{1}{8}\right) \tau^2 \\
&\stackrel{(iii)}{=} 1 - \frac{\mu}{4L^{\frac{3}{2}}} + \left(\frac{2}{\sqrt{L}\mu} + \frac{1}{8}\right) \tau^2 \\
&\stackrel{(iv)}{\leq} 1 - \frac{\mu}{8L^{\frac{3}{2}}} = \alpha, \tag{20}
\end{aligned}$$

where in (i) we plugged in the choice of η ; (ii) is because we have assumed that $L > 1$; (iii) plugged in the choice of stepsize; (iv) uses the definition of τ in (17); and the last equality comes from the definition of α in (17). Plugging (20) to (19), we obtain $f(\theta^{t+1}) \leq (\alpha)^{t+1} f(\theta^0)$.

Second, we show that (15) holds for $t + 1$. We have:

$$\begin{aligned}
\|g^{t+1} - q^t\|_\infty &\leq \|g^{t+1} - g^t\| + \|g^t - q^t\| \\
&\stackrel{(i)}{\leq} L\eta\|q^t\| + \tau\gamma^{t-1} \\
&\stackrel{(ii)}{\leq} L\eta(\|g^t\| + D\|g^t - q^t\|_\infty) + \tau\gamma^{t-1} \\
&\stackrel{(iii)}{\leq} L\eta\sqrt{2L}\sqrt{f(\theta^t)} + (1 + DL\eta)\tau\gamma^{t-1} \\
&\stackrel{(iv)}{\leq} L\eta\sqrt{L}\sqrt{(\alpha)^t f(\theta^0)} + (1 + DL\eta)\tau\gamma^{t-1} \\
&\stackrel{(v)}{=} (L\eta\sqrt{2L} + (1 + DL\eta)\tau)\gamma^{t-1} \\
&= \left(\frac{\sqrt{2}}{2} + (1 + \frac{D}{2})\tau\right)\gamma^{t-1} \\
&\stackrel{(vi)}{\leq} \left(\frac{\sqrt{2}}{2} + \frac{3}{2}D\tau\right)\gamma^{t-1} \stackrel{(vii)}{\leq} \frac{9}{10}\gamma^{t-1} \stackrel{(viii)}{\leq} \sqrt{\alpha}\gamma^t = \gamma^{t+1}
\end{aligned}$$

where (i) uses Assumption 2 and the induction assumption; (ii) uses the triangle inequality to decompose $\|q^t\|$ and uses the relationship between ℓ_∞ and ℓ_2 ; (iii) is from Assumption 1 and induction assumption (15); (iv) uses the induction assumption (14); (v) uses the definition of γ^t ; (vi) is because $\mu < 1$; (vii) comes from the fact that $\tau < \frac{1}{10D}$ in (17); (viii) compares the choice of α in (17). Thus, we obtain $\|g^{t+1} - q^t\|_\infty \leq \sqrt{(\alpha)^{t+1} f(\theta^0)} = \gamma^t$. Then by Lemma 4.1, we have (15) holds for $t + 1$.

Now we have proved by induction that (14) and (15) hold. So for $t > 0$, there is

$$f(\theta^t) \leq (\alpha)^t f(\theta^0), \text{ where } \alpha = 1 - \frac{\mu}{8L^{\frac{3}{2}}}.$$

Thus, to compute an ϵ -optimal solution, the total number of iterations required is $\log(f(\theta^0)/\epsilon)/\log(1 - \mu/8L^{\frac{3}{2}})$. Since in each iteration, each agent k transmits a length- D vector q_k^t , it follows that the total number of bits each agent needs to communicate is $D \log(f(\theta^0)/\epsilon)/\log(1 - \mu/8L^{\frac{3}{2}})$ bits. Notice that $\log(1/(1 - \mu/8L^{\frac{3}{2}})) = -\log(1 - \mu/8L^{\frac{3}{2}}) \sim 8L^{\frac{3}{2}}/\mu$, so we can derive the simplified total number of bits as $bD \cdot \frac{8L^{\frac{3}{2}}}{\mu} \log(f(\theta^0)/\epsilon)$.

D The Proof of Theorem 4.2

First, let us provide the sketch of the proof. Denote

$$g^t := \sum_{k=1}^K g_k^t, \quad q^t := \sum_{k=1}^K q_k^t, \quad \tilde{g}^t := B^\top g^t.$$

Using the above notation, the agents' local update step (i.e., the 'Update' step in Alg. 1) can be expressed as:

$$\theta_k^{t+1} = \theta_k^t - \eta B^\top q^t, \forall k. \quad (21)$$

Step 1: We show that the loss function decreases linearly if all the agents update parameters using the direction \tilde{g}^t , as follows:

$$f(\theta^t - \eta\tilde{g}^t) \leq \left(1 - \frac{1}{\kappa^2}\right) f(\theta^t).$$

Step 2: Let $\tau = \frac{1}{2^b - 1}$, we show by induction that for $t = 1, \dots$, the following inequalities hold true:

$$(1) f(\theta^t) \leq (\alpha)^t f(\theta^0), \text{ for some } 0 < \alpha < 1. \quad (22)$$

$$(2) \|g^t - q^t\|_\infty \leq \tau\gamma^{t-1}, \|\tilde{g}^t - B^\top q^t\|_\infty \leq \tau H\|B^\top\|_\infty\gamma^{t-1}, \text{ for some } \gamma^{t-1} > 0, \quad (23)$$

where $\{\gamma^t\}_{t=1}^\infty$ is a sequence of positive numbers.

We prove **Step 1** first. At t -the iteration, Let us first expand the objective function as:

$$f(\theta^t - \eta\tilde{g}^t) = \frac{1}{2}\|A(\theta^t - \eta\tilde{g}^t) - b\|^2 = f(\theta^t) - \langle A\theta^t - b, \eta A\tilde{g}^t \rangle + \frac{1}{2}\eta^2\|A\tilde{g}^t\|^2. \quad (24)$$

To proceed, let us provide explicit expressions for g^t and \tilde{g}^t :

$$g^t = B\nabla f(\theta^t) = BA^\top(A\theta^t - b) \quad (25)$$

$$\tilde{g}^t = B^\top B\nabla f(\theta^t) = B^\top BA^\top(A\theta^t - b). \quad (26)$$

By using the above, the inner product in (24) can be bounded as follows

$$-\langle A\theta^t - b, \eta A\tilde{g}^t \rangle = -\eta\langle A\theta^t - b, AB^\top BA^\top(A\theta^t - b) \rangle \leq -\eta\sigma_{\min}(Z)\|A\theta^t - b\|^2. \quad (27)$$

Using the above relations, we can further bound the descent of the objective function as

$$\begin{aligned} f(\theta^t - \eta\tilde{g}^t) &\stackrel{(i)}{\leq} f(\theta^t) - \eta\sigma_{\min}(Z)\|A\theta^t - b\|^2 + \frac{1}{2}\eta^2\|AB^\top BA^\top(A\theta^t - b)\|^2 \\ &\stackrel{(ii)}{\leq} f(\theta^t) - \eta\sigma_{\min}(Z)\|A\theta^t - b\|^2 + \frac{1}{2}\eta^2\sigma_{\max}^2(Z)\|A\theta^t - b\|^2 \\ &\stackrel{(iii)}{=} f(\theta^t) - 2\eta\sigma_{\min}(Z)f(\theta^t) + \eta^2\sigma_{\max}^2(Z)f(\theta^t) \\ &= (1 - \eta\sigma_{\min}(Z))f(\theta^t) \\ &= \left(1 - \frac{1}{\kappa^2}\right)f(\theta^t), \end{aligned} \quad (28)$$

where (i) comes from plugging (27) and (26) into (24); (ii) extracts the largest eigenvalue of Z ; (iii) is due to the definition of the objective function; the last two equalities hold due to the definition $\eta = \frac{\sigma_{\min}(Z)}{\sigma_{\max}^2(Z)}$.

Next, we prove **Step 2**. To begin with, let us define:

$$\alpha := 1 - \frac{1}{2\kappa^4}, \quad \lambda := 6\sqrt{2} \cdot \frac{\sigma_{\min}(Z)}{\sqrt{\sigma_{\max}(Z)}}, \quad \gamma^t := \lambda\sqrt{(\alpha)^{t+1}f(\theta^0)}. \quad (29)$$

Further let us set

$$b = \max\left(\log_2\left(\frac{1}{\tau} + 1\right), b_0\right), \quad C = \max\left(\frac{1}{\lambda^2\tau^2D^2}, \frac{12}{\lambda}\right) \quad \text{with } \tau \text{ given by} \quad (30)$$

$$\tau := \min\left(\frac{1}{\sqrt{2\kappa^4\lambda^2D^2H^2\|B^\top\|_\infty^2s_{\max}^2(A)\left(\frac{\sigma_{\min}^2(Z)}{2\sigma_{\max}^4(Z)} + \frac{1}{\sigma_{\max}^2(Z)}\right)}}, \frac{1}{6\left(1 + \frac{H}{\kappa}\right)}\right).$$

First, we analyze (22) for $t = 1$. We have the following relations:

$$\begin{aligned} f(\theta^1) &= f(\theta^0 - \eta B^\top q^0) - f(\theta^0 - \eta\tilde{g}^0) + f(\theta^0 - \eta\tilde{g}^0) \\ &\stackrel{(i)}{\leq} \frac{1}{2}\|\eta A(B^\top q^0 - \tilde{g}^0)\|^2 + \langle \eta A(B^\top q^0 - \tilde{g}^0), A(\theta^0 - \eta\tilde{g}^0) - b \rangle + (1 - \eta\sigma_{\min}(Z))f(\theta^0) \\ &\stackrel{(ii)}{\leq} \frac{1}{2}\|\eta A(B^\top q^0 - \tilde{g}^0)\|^2 + \eta\left(\beta\|A(B^\top q^0 - \tilde{g}^0)\|^2 + \frac{1}{2\beta}\|A(\theta^0 - \eta\tilde{g}^0) - b\|^2\right) + (1 - \eta\sigma_{\min}(Z))f(\theta^0) \\ &\stackrel{(iii)}{\leq} \left(\frac{1}{2}\eta^2 + \beta\eta\right)s_{\max}^2(A)\|B^\top q^0 - \tilde{g}^0\|^2 + \frac{\eta}{\beta}(1 - \eta\sigma_{\min}(Z))f(\theta^0) + (1 - \eta\sigma_{\min}(Z))f(\theta^0) \\ &= \left(\frac{1}{2}\eta^2 + \beta\eta\right)s_{\max}^2(A)\|B^\top q^0 - \tilde{g}^0\|^2 + \left(1 + \frac{\eta}{\beta}\right)(1 - \eta\sigma_{\min}(Z))f(\theta^0) \\ &\stackrel{(iv)}{=} \left(\frac{1}{2}\eta^2 + \frac{\eta}{\sigma_{\min}(Z)}\right)s_{\max}^2(A)\|B^\top q^t - \tilde{g}^0\|^2 + (1 - \eta^2\sigma_{\min}^2(Z))f(\theta^0) \end{aligned}$$

$$\begin{aligned}
&\stackrel{(v)}{\leq} D^2 \left(\frac{1}{2} \eta^2 + \frac{\eta}{\sigma_{\min}(Z)} \right) s_{\max}^2(A) \|B^\top q^0 - \tilde{g}^0\|_\infty^2 + (1 - \eta^2 \sigma_{\min}^2(Z)) f(\theta^0) \\
&\stackrel{(vi)}{\leq} D^2 H^2 \|B^\top\|_\infty^2 \left(\frac{1}{2} \eta^2 + \frac{\eta}{\sigma_{\min}(Z)} \right) s_{\max}^2(A) \|q^0 - g^0\|_\infty^2 + (1 - \eta^2 \sigma_{\min}^2(Z)) f(\theta^0) \\
&\stackrel{(vii)}{\leq} K D^2 H^2 \|B^\top\|_\infty^2 \left(\frac{1}{2} \eta^2 + \frac{\eta}{\sigma_{\min}(Z)} \right) s_{\max}^2(A) \sum_{k=1}^K \|q_k^0 - g_k^0\|_\infty^2 + (1 - \eta^2 \sigma_{\min}^2(Z)) f(\theta^0) \\
&\stackrel{(viii)}{\leq} H^2 \|B^\top\|_\infty^2 \left(\frac{1}{2} \eta^2 + \frac{\eta}{\sigma_{\min}(Z)} \right) s_{\max}^2(A) \frac{f(\theta^0)}{C^2 D^4} + (1 - \eta^2 \sigma_{\min}^2(Z)) f(\theta^0) \\
&= \left(1 - \eta^2 \sigma_{\min}^2(Z) + H^2 \|B^\top\|_\infty^2 \left(\frac{1}{2} \eta^2 + \frac{\eta}{\sigma_{\min}(Z)} \right) \frac{s_{\max}^2(A)}{C^2 D^4} \right) f(\theta^0) \\
&\stackrel{(ix)}{=} \left(1 - \frac{1}{\kappa^4} + H^2 \|B^\top\|_\infty^2 \left(\frac{\sigma_{\min}^2(Z)}{2\sigma_{\max}^4(Z)} + \frac{1}{\sigma_{\max}^2(Z)} \right) s_{\max}^2(A) \lambda^2 \tau^2 / D^2 \right) f(\theta^0) \\
&\stackrel{(x)}{\leq} \left(1 - \frac{1}{2\kappa^4} \right) f(\theta^0) = \alpha f(\theta^0),
\end{aligned}$$

where (i) explicitly expands the $f(\cdot)$ function, and uses (28); (ii) applies the Young inequality with constant β ; (iii) extracts the largest singular value of A , and uses (28); (iv) set $\beta = \frac{1}{\sigma_{\min}(Z)}$; (v) uses the relationship between ℓ_2 and ℓ_∞ norm; (vi) uses the fact that $\|B^\top q^0 - \tilde{g}^0\| \leq H^2 \|\tilde{B}\|_\infty \|q^0 - g^0\|_\infty$; (vii) uses the Cauchy-Schwartz inequality; (viii) uses the condition $\|q_k^0 - g_k^0\|_\infty \leq \frac{\sqrt{f(\theta^0)}}{CD\sqrt{K}}$ in Algorithm 1; (vii) uses the induction assumption (22) and the definition of γ^{t-1} in (29); (viii) plug in the choice of stepsize; (ix) comes from the choice of λ and τ . So we have showed (22) holds for $t = 1$; (ix) plug in the choice of stepsize and constant $C \geq \frac{1}{\lambda^2 \tau^2}$; (x) uses the choice of λ and τ .

Next, let us analyze (23) for $t = 1$. The idea is that, if we can show that $\|g^1 - q^0\|_\infty \leq \gamma^0$, then we will be able to use Lemma 4.1 to show (23). More specifically, we have:

$$\begin{aligned}
&\|g^1 - g^1\| \stackrel{(i)}{=} \|\text{quant}(g^1, q^0, \gamma^0, b) - g^1\|_\infty \stackrel{(ii)}{\leq} \tau \gamma^0, \\
&\|B^\top q^1 - \tilde{g}^1\|_\infty \stackrel{(iii)}{=} \|B^\top \text{quant}(g^1, q^0, \gamma^0, b) - B^\top g^1\|_\infty \stackrel{(iv)}{\leq} \tau H \|B^\top\|_\infty \gamma^0,
\end{aligned}$$

where (i) and (iii) come from the ‘Quantize’ step in Algorithm 1; (ii) and (iv) are from the two inequalities in Lemma 4.1 (assuming that $\|g^1 - q^0\|_\infty \leq \gamma^0$ holds).

Next, we show $\|g^1 - q^0\|_\infty \leq \gamma^0$. We observe that:

$$\begin{aligned}
\|g^1 - q^0\|_\infty &\stackrel{(i)}{\leq} \|g^1 - g^0\|_\infty + \|g^0 - q^0\|_\infty \\
&\stackrel{(ii)}{\leq} \eta \|BA^\top AB^\top q^0\| + \|g^0 - q^0\|_\infty \\
&\stackrel{(iii)}{\leq} \eta \sigma_{\max}(Z) \|q^0\| + \|g^0 - q^0\|_\infty \\
&\stackrel{(iv)}{\leq} \eta \sigma_{\max}(Z) (\|g^0\| + \|q^0 - g^0\|) + \|g^0 - q^0\|_\infty \\
&\stackrel{(v)}{\leq} \eta \sigma_{\max}(Z) (\|BA^\top(A\theta^0 - b)\| + H \|q^0 - g^0\|_\infty) + \|g^0 - q^0\|_\infty \\
&\stackrel{(vi)}{\leq} \eta \sigma_{\max}^{\frac{3}{2}}(Z) \|A\theta^0 - b\| + (1 + \eta H \sigma_{\max}(Z)) \|g^0 - q^0\|_\infty \\
&\stackrel{(vii)}{\leq} \eta \sigma_{\max}^{\frac{3}{2}}(Z) \|A\theta^0 - b\| + (1 + \eta H \sigma_{\max}(Z)) \sum_{k=1}^K \|g_k^0 - q_k^0\|_\infty \\
&\stackrel{(viii)}{\leq} \eta \sigma_{\max}^{\frac{3}{2}}(Z) \sqrt{2f(\theta^0)} + (1 + \eta H \sigma_{\max}(Z)) \sum_{k=1}^K \frac{\sqrt{f_k(\theta^0)}}{CD^3 \sqrt{K}} \\
&\stackrel{(ix)}{\leq} \eta \sigma_{\max}^{\frac{3}{2}}(Z) \sqrt{2f(\theta^0)} + (1 + \eta H \sigma_{\max}(Z)) \frac{\sqrt{f(\theta^0)}}{CD^3}
\end{aligned}$$

$$\begin{aligned}
&\stackrel{(x)}{=} \sqrt{f(\theta^0)} \left(\frac{1}{CD^3} + \frac{H}{CD^3\kappa} + \frac{\sqrt{2}\sigma_{\min}(Z)}{\sqrt{\sigma_{\max}(Z)}} \right) \\
&\stackrel{(xi)}{\leq} \lambda \sqrt{f(\theta^0)} \left(\frac{2}{\lambda C} + \frac{1}{6} \right) \stackrel{(xii)}{\leq} \frac{1}{3} \lambda \sqrt{f(\theta^0)} \leq \lambda \sqrt{\alpha f(\theta^0)} = \gamma^0
\end{aligned}$$

where (i) is due to the triangle inequality; (ii) expands the expression of g^1 and g^0 in (25), uses the relation between ℓ_2 norm and ℓ_∞ norm and uses the update rule (21); (iii) uses the fact that non-zero eigen values of $BA^\top AB^\top$ and Z are the same and extracts the largest eigen value of Z ; (iv) uses triangle inequality; (v) uses the relationship between ℓ_2 and ℓ_∞ norm; (vi) is because $s_{\max}^2(BA^\top) = \sigma_{\max}(Z)$ and extract the largest singular value of BA^\top ; (vii) uses triangle inequality; (viii) uses the initial condition $\|q_k^0 - g_k^0\|_\infty \leq \frac{\sqrt{f_k(\theta^0)}}{C \cdot D \sqrt{K}}$ in Algorithm 1; (ix) uses Cauchy-Schwartz inequality; (x) plug in the choice of stepsize; (xi) is because $\kappa \geq 1, H \leq D, D \geq 1$ and the choice of λ ; (xii) is from $C \geq \frac{12}{\lambda}$; the last inequality comes from $\sqrt{\alpha} \geq \frac{\sqrt{2}}{2} > \frac{1}{3}$ since $\kappa \geq 1$. So we can show $\|g^1 - q^0\|_\infty \leq \gamma^0$.

Next, we will show (22) holds for $t + 1$ by induction, based on the base assumption that (22) and (23) holds for t . We have the following series of relations:

$$\begin{aligned}
f(\theta^{t+1}) &= f(\theta^t - \eta B^\top q^t) - f(\theta^t - \eta \tilde{g}^t) + f(\theta^t - \eta \tilde{g}^t) \\
&\stackrel{(i)}{\leq} \frac{1}{2} \|\eta A(B^\top q^t - \tilde{g}^t)\|^2 + \langle \eta A(B^\top q^t - \tilde{g}^t), A(\theta^t - \eta \tilde{g}^t) - b \rangle + (1 - \eta \sigma_{\min}(Z)) f(\theta^t) \\
&\stackrel{(ii)}{\leq} \frac{1}{2} \|\eta A(B^\top q^t - \tilde{g}^t)\|^2 + \eta \left(\beta \|A(B^\top q^t - \tilde{g}^t)\|^2 + \frac{1}{2\beta} \|A(\theta^t - \eta \tilde{g}^t) - b\|^2 \right) + (1 - \eta \sigma_{\min}(Z)) f(\theta^t) \\
&\stackrel{(iii)}{\leq} \left(\frac{1}{2} \eta^2 + \beta \eta \right) s_{\max}^2(A) \|B^\top q^t - \tilde{g}^t\|^2 + \frac{\eta}{\beta} (1 - \eta \sigma_{\min}(Z)) f(\theta^t) + (1 - \eta \sigma_{\min}(Z)) f(\theta^t) \\
&= \left(\frac{1}{2} \eta^2 + \beta \eta \right) s_{\max}^2(A) \|B^\top q^t - \tilde{g}^t\|^2 + \left(1 + \frac{\eta}{\beta} \right) (1 - \eta \sigma_{\min}(Z)) f(\theta^t) \\
&\stackrel{(iv)}{=} \left(\frac{1}{2} \eta^2 + \frac{\eta}{\sigma_{\min}(Z)} \right) s_{\max}^2(A) \|B^\top q^t - \tilde{g}^t\|^2 + (1 - \eta^2 \sigma_{\min}^2(Z)) f(\theta^t) \\
&\stackrel{(v)}{\leq} D^2 \left(\frac{1}{2} \eta^2 + \frac{\eta}{\sigma_{\min}(Z)} \right) s_{\max}^2(A) \|B^\top q^t - \tilde{g}^t\|_\infty^2 + (1 - \eta^2 \sigma_{\min}^2(Z)) f(\theta^t) \\
&\stackrel{(vi)}{\leq} D^2 H^2 \|B^\top\|_\infty^2 \left(\frac{1}{2} \eta^2 + \frac{\eta}{\sigma_{\min}(Z)} \right) s_{\max}^2(A) (\tau \gamma^{t-1})^2 + (1 - \eta^2 \sigma_{\min}^2(Z)) f(\theta^t) \\
&\stackrel{(vii)}{\leq} \left(1 - \eta^2 \sigma_{\min}^2(Z) + D^2 H^2 \|B^\top\|_\infty^2 \left(\frac{1}{2} \eta^2 + \frac{\eta}{\sigma_{\min}(Z)} \right) s_{\max}^2(A) \lambda^2 \tau^2 \right) (\alpha)^t f(\theta^0) \\
&\stackrel{(viii)}{=} \left(1 - \frac{1}{\kappa^4} + D^2 H^2 \|B^\top\|_\infty^2 \left(\frac{\sigma_{\min}^2(Z)}{2\sigma_{\max}^4(Z)} + \frac{1}{\sigma_{\max}^2(Z)} \right) s_{\max}^2(A) \lambda^2 \tau^2 \right) (\alpha)^t f(\theta^0) \\
&\stackrel{(ix)}{\leq} \left(1 - \frac{1}{2\kappa^4} \right) (\alpha)^t f(\theta^0) = (\alpha)^{t+1} f(\theta^0),
\end{aligned}$$

where (i) we have explicitly expands the $f(\cdot)$ function, and have used (28); (ii) applies the Young inequality with constant β ; (iii) extracts the largest singular value of A , and uses (28); (iv) set $\beta = \frac{1}{\sigma_{\min}(Z)}$; (v) uses the relationship between ℓ_2 and ℓ_∞ norm; (vi) uses the second inequality in induction assumption (23); (vii) uses the induction assumption (22) and the definition of γ^{t-1} in (29); (viii) plug in the choice of stepsize; (ix) comes from the choice of λ and τ .

Finally, we will show (23) holds for $t + 1$ by induction, again base assumptions (22) holds for t and (23) holds for t . We have:

$$\begin{aligned}
\|g^{t+1} - q^t\|_\infty &\stackrel{(i)}{\leq} \|g^{t+1} - g^t\|_\infty + \|g^t - q^t\|_\infty \\
&\stackrel{(ii)}{\leq} \|BA^\top A(\theta^{t+1} - \theta^t)\| + \tau \gamma^{t-1} \\
&\stackrel{(iii)}{=} \eta \|BA^\top AB^\top q^t\| + \tau \gamma^{t-1}
\end{aligned}$$

$$\begin{aligned}
&\stackrel{(iv)}{\leq} \eta\sigma_{\max}(Z)\|q^t\| + \tau\gamma^{t-1} \\
&\stackrel{(v)}{\leq} \eta\sigma_{\max}(Z)(\|q^t - g^t\| + \|g^t\|) + \tau\gamma^{t-1} \\
&\stackrel{(vi)}{\leq} \eta\sigma_{\max}(Z)(H\|q^t - g^t\|_{\infty} + \|g^t\|) + \tau\gamma^{t-1} \\
&\stackrel{(vii)}{\leq} \tau\gamma^{t-1}(1 + \eta H\sigma_{\max}(Z)) + \eta\sigma_{\max}(Z)\|g^t\| \\
&= \tau\gamma^{t-1}(1 + \eta H\sigma_{\max}(Z)) + \eta\sigma_{\max}(Z)\|BA^{\top}(A\theta^t - b)\| \\
&\stackrel{(viii)}{\leq} \tau\gamma^{t-1}(1 + \eta H\sigma_{\max}(Z)) + \eta\sigma_{\max}^{\frac{3}{2}}(Z)\sqrt{2f(\theta^t)} \\
&\stackrel{(ix)}{=} \tau\gamma^{t-1}\left(1 + \frac{H}{\kappa}\right) + \frac{\sqrt{2}\sigma_{\min}(Z)}{\sqrt{\sigma_{\max}(Z)}}\sqrt{f(\theta^t)} \\
&\stackrel{(x)}{\leq} \left(\tau\lambda\left(1 + \frac{H}{\kappa}\right) + \frac{\sqrt{2}\sigma_{\min}(Z)}{\sqrt{\sigma_{\max}(Z)}}\right)\sqrt{(\alpha)^t f(\theta^0)} \\
&= \lambda\left(\tau\left(1 + \frac{H}{\kappa}\right) + \frac{\sqrt{2}\sigma_{\min}(Z)}{\lambda\sqrt{\sigma_{\max}(Z)}}\right)\sqrt{(\alpha)^t f(\theta^0)} \\
&\stackrel{(xi)}{\leq} \frac{1}{3}\lambda\sqrt{(\alpha)^t f(\theta^0)} \stackrel{(xii)}{\leq} \lambda\sqrt{(\alpha)^{t+1} f(\theta^0)} = \gamma^t,
\end{aligned}$$

where (i) uses the triangle inequality; (ii) expands the first term by (25), uses the relationship between ℓ_2 and ℓ_{∞} norm, and uses the first inequality in induction assumption (23); (iii) plug in the update of parameter: $\theta^{t+1} = \theta^t - \eta B^{\top} q^t$; (iv) comes from the fact that non-zero singular values of Z and $BA^{\top}AB^{\top}$ are the same and extracts the largest singular value of Z ; (v) uses the triangle inequality; (vi) is due to the relationship between ℓ_2 and ℓ_{∞} norm; (vii) uses the first inequality in induction assumption (23); (viii) uses the fact that $s_{\max}^2(BA^{\top}) = \sigma_{\max}(Z)$ and extracts $s_{\max}(BA^{\top})$; (ix) plug in the choice of stepsize; (x) uses the definition of γ^{t-1} in (29) and induction assumption in (22); (xi) is because the choice of τ and λ ; (xii) is because $\sqrt{\alpha} = \sqrt{1 - \frac{1}{2\kappa^2}} \geq \frac{\sqrt{2}}{2} > \frac{1}{3}$ since $\kappa \geq 1$. Thus, we obtain $\|g^{t+1} - q^t\|_{\infty} \leq \lambda\sqrt{(\alpha)^{t+1} f(\theta^0)} = \gamma^t$. Then by Lemma 4.1, with the correspondence that $c = g^{t+1}$, $p = q^t$, $r = \gamma^t$, we can obtain

$$\|q^{t+1} - g^{t+1}\| \stackrel{(i)}{=} \|\text{quant}(g^{t+1}, q^t, \gamma^t, b) - g^{t+1}\|_{\infty} \stackrel{(ii)}{\leq} \tau\gamma^t,$$

$$\|B^{\top}q^{t+1} - \tilde{g}^{t+1}\|_{\infty} \stackrel{(iii)}{=} \|B^{\top}\text{quant}(g^{t+1}, q^t, \gamma^0, b) - B^{\top}g^{t+1}\|_{\infty} \stackrel{(iv)}{\leq} \tau D\|B^{\top}\|_{\infty}\gamma^t,$$

where (i) and (iii) come from the ‘Quantize’ step in Algorithm 1; (ii) and (iv) are from the two relations in Lemma 4.1 (since we have proved that $\|g^{t+1} - q^t\|_{\infty} \leq \gamma^t$ holds).

Now we have proved that (14) and (15) hold. So for $t > 0$, there is

$$f(\theta^t) \leq (\alpha)^t f(\theta^0), \quad \text{where} \quad \alpha = 1 - \frac{1}{2\kappa^4}.$$

Thus, if we want the objective function to compute an ϵ -optimal solution, the total number of iterations is $\log(f(\theta^0)/\epsilon)/\log(1/(1 - \frac{1}{2\kappa^4}))$. Since in each iteration, each agent k transmits a length- H vector q_k^t , so we conclude that the total number of bits each node needs to communicate is $\log(f(\theta^0)/\epsilon)/\log(1/(1 - \frac{1}{2\kappa^4}))$ bits. Notice that $\log(1/(1 - \frac{1}{2\kappa^4})) = -\log(1 - \frac{1}{2\kappa^4}) \sim 2\kappa^4$, so we can derive the simplified total number of bits as $2\kappa^4 \cdot \log(f(\theta^0)/\epsilon)$.

E Proof for Proposition 1

Now we prove Proposition 1. We first state two lemmas that will be used.

Lemma E.1. *Rudelson and Vershynin [2010]* Let X be a $H \times N$ matrix whose entries are independent standard normal random variables. Then

$$\mathbb{P}\left(\sqrt{H} - \sqrt{N} - t \leq s_{\min}(X) \leq s_{\max}(X) \leq \sqrt{H} + \sqrt{N} + t\right) \geq 1 - 2e^{-t^2/2}, \quad t \geq 0.$$

Lemma E.2. Suppose $A_i^\top \in \mathbb{R}^D$ follows $N(\mu, \Sigma)$. If Σ is a diagonal matrix, then each element of A_i is independent.

We first write down the SVD decomposition of A^\top :

$$A^\top = V\Sigma W,$$

where $V \in \mathbb{R}^{D \times N}$, $\Sigma \in \mathbb{R}^{N \times N}$, $W \in \mathbb{R}^{N \times N}$. Since A is full rank, denote $\hat{\kappa} := \frac{s_{\max}(A^\top)}{s_{\min}(A^\top)} > 0$. Consider each entry in $BA^\top = BV\Sigma W$. It is clear that each entry of BV follows normal distribution because $(BV)_{hi}$ is linear combination of variables that are standard normal. Then we aim to show that each entry in BV follows standard normal distribution. Notice that $\text{vec}(BV) = (B_1V, B_2V, \dots, B_HV)^\top \in \mathbb{R}^{HN}$, where B_h is each row of B . Denote $U_{ij} = \text{Cov}(B_iV, B_jV) \in \mathbb{R}^{N \times N}$, we can obtain that

$$\begin{aligned} \mathbb{E}[\text{vec}(BV)] &= (\mathbb{E}[B_1]V; \mathbb{E}[B_2]V; \dots, \mathbb{E}[B_H]V)^\top = \mathbf{0}, \\ \text{Cov}(\text{vec}(BV)) &= \begin{pmatrix} U_{11} & U_{12} & U_{13} & \cdots & U_{1N} \\ U_{21} & U_{22} & U_{23} & \cdots & U_{2N} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ U_{N1} & U_{N2} & U_{N3} & \cdots & U_{NN} \end{pmatrix} \end{aligned}$$

Notice that $U_{ij} = \text{Cov}(X^\top B_i^\top, X^\top B_j^\top) = X^\top \text{Cov}(B_i^\top, B_j^\top) X$. Since B_i^\top and B_j^\top are independent, so $U_{ij} = \mathbf{0}$, $i \neq j$. Since $X^\top X = \mathbf{I}_N$, we have $U_{ij} = \text{Cov}(B_i^\top, B_i^\top) = \mathbf{I}_N$ for $i = 1, 2, \dots, N$. Then we obtain $\text{Cov}(\text{vec}(BV)) = \mathbf{I}_{HN}$. From Lemma E.2, each entry in BV are independent. By Lemma E.1, we know the condition number of BV , which is $\kappa(BV)$ is independent of D , and we have

$$\mathbb{P}\left(\kappa(BV) \leq \frac{\sqrt{H} + \sqrt{N} + t}{\sqrt{H} - \sqrt{N} - t}\right) \geq 1 - 2e^{-t^2/2}, \quad t \geq 0.$$

Now we consider $Z = AB^\top BA^\top$. Notice the non-zero eigen values of Z are the same as $BA^\top AB^\top$, then we consider $BA^\top AB^\top$,

$$BA^\top AB^\top = BV\Sigma W W^\top \Sigma^\top V^\top B^\top = BV\Sigma^2 V^\top B^\top.$$

Notice the non-zero eigen values of $BV\Sigma^2 V^\top B^\top$ are the same as $\Sigma^2 B V V^\top B^\top$, we can derive

$$\kappa(Z) = \kappa(\Sigma^2 B V V^\top B^\top) \leq \hat{\kappa}^2 \cdot \kappa^2(BV),$$

where $\hat{\kappa}$ is condition number of A and the inequality is because the property of square and invertible matrix. So we have with probability at least $1 - 2e^{-t^2/2}$,

$$\kappa(Z) \leq \hat{\kappa}^2 \left(\frac{\sqrt{H} + \sqrt{N} + t}{\sqrt{H} - \sqrt{N} - t} \right)^2.$$

F The Proof of Theorem 4.3 and Details about Algorithm 2

F.1 The CHOCO-GOSSIP Protocol

For completeness, we describe the CHOCO-GOSSIP protocol [Koloskova et al., 2019] for decentralized average consensus with compressed communication as follows. Notice the number of gossip rounds T_g is dependent to the error ϵ as discussed in Fact 1. Note that the protocol uses the compressor $Q(\cdot)$ for compressed communication, for example, this can be the randomized quantizer, see [Koloskova et al., 2019] for other examples. The protocol is summarized below:

Algorithm 3 CHOCO-GOSSIP

Input: step size γ ; initial vectors g_1^0, \dots, g_K^0 ; gossip rounds T_g ; compressor $Q(\cdot)$; mixing matrix W ; neighbor sets $\mathcal{N}_1, \dots, \mathcal{N}_K$.

Initialize: $\hat{g}_i^0 = \mathbf{0}$, $\forall i \in [K]$.

for t **in** $0, \dots, T_g - 1$ **do**

Compress: Each agent i compress the difference, $q_i^t = Q(\hat{g}_i^t - g_i^t)$

Communicate: Each agent i receives q_j^t from neighbor $j \in \mathcal{N}_i$ and update $\hat{g}_j^{t+1} = \hat{g}_j^t + q_j^t$

Aggregation: Each agent i combines received vectors, i.e.,

$$g_i^{t+1} = g_i^t + \gamma \sum_{j \in \mathcal{N}_i} w_{ij} (\hat{g}_j^{t+1} - \hat{g}_i^{t+1})$$

end for

Output: $g_i^{T_g}$ at each agent i .

To derive the number of bits required in Theorem 4.3, we focus on using the random quantizer for $Q(\cdot)$, i.e., for $x \in \mathbb{R}^d$, $s \in \mathbb{N}_+$ and $\tau = (1 + \min\{d/s^2, \sqrt{d}/s\})$, we have

$$Q(x) = \frac{\text{sign}(x) \cdot \|x\|}{s\tau} \cdot \left\lfloor s \frac{|x|}{\|x\|} + \xi \right\rfloor,$$

where $\xi \sim \text{Uniform}[0, 1]^d$ and sending $Q(x)$ across the network requires $d \log(s + 1) + d + 64$ bits of communication.

F.2 The Proof of Theorem 4.3

Since $F(x) = Bx$, $\tilde{F}(y) = B^\top y$, the consensus error can be expressed as $E_i^t := \bar{g}_i^t - \frac{1}{n} \sum_{j=1}^n B \nabla f_j(\theta_j^t)$. We also denote the deviation of locally computed gradient as $\Delta G_i^t := \frac{1}{n} \sum_{j=1}^n B (\nabla f_j(\theta_j^t) - \nabla f_j(\theta_i^t))$. Note that $\|E_i^t\| \leq \bar{\epsilon}/(t + 1)$. We first observe the updated iterate,

$$\begin{aligned} \theta_i^{t+1} &= \theta_i^t - \eta B^\top \bar{g}_i^t - \frac{1}{n} \sum_{j=1}^n B^\top B \nabla f_j(\theta_j^t) + \frac{1}{n} \sum_{j=1}^n B^\top B \nabla f_j(\theta_j^t) \\ &= \theta_i^t - \eta \left[B^\top E_i^t + \frac{1}{n} B^\top B \sum_{j=1}^n (\nabla f_j(\theta_j^t) + \nabla f_j(\theta_i^t) - \nabla f_j(\theta_i^t)) \right] \\ &= \theta_i^t - \eta [B^\top E_i^t + B^\top \Delta G_i^t + B^\top B \nabla f(\theta_i^t)] \end{aligned}$$

Next, we observe that the objective function value evolves as,

$$\begin{aligned} f(\theta_i^{t+1}) &= \frac{1}{2} \sum_{j=1}^n \|A_j \theta_i^{t+1} - b_j\|^2 \\ &= \frac{1}{2} \sum_{j=1}^n \left\{ \|A_j \theta_i^t - b_j\|^2 - 2\eta \langle A_j^\top (A_j \theta_i^t - b_j) \mid B^\top E_i^t + B^\top \Delta G_i^t + B^\top B \nabla f(\theta_i^t) \rangle \right\} \end{aligned}$$

$$\begin{aligned}
& + \eta^2 \|A_j [B^\top E_i^t + B^\top \Delta G_i^t + B^\top B \nabla f(\theta_i^t)]\|^2 \Big\} \\
& \leq f(\theta_i^t) - \eta \langle \nabla f(\theta_i^t) \mid B^\top B \nabla f(\theta_i^t) \rangle + \frac{1}{2} \sum_{j=1}^n 2\eta^2 \|A_j B^\top B \nabla f(\theta_i^t)\|^2 \\
& \quad + \frac{1}{2} \sum_{j=1}^n \left\{ -2\eta \langle A_j \theta_i^t - \mathbf{b}_j \mid A_j B^\top (E_i^t + \Delta G_i^t) \rangle + 2\eta^2 \|A_j B^\top (E_i^t + \Delta G_i^t)\|^2 \right\} \\
& \leq (1 - 2\sigma_{\min}(Z)\eta + 2\sigma_{\max}^2(Z)\eta^2) f(\theta_i^t) \\
& \quad + \sum_{j=1}^n \left\{ -\eta \langle A_j \theta_i^t - \mathbf{b}_j \mid A_j B^\top (E_i^t + \Delta G_i^t) \rangle + \eta^2 \|A_j B^\top (E_i^t + \Delta G_i^t)\|^2 \right\} \quad (31)
\end{aligned}$$

where the last inequality requires the observations

$$\begin{aligned}
\eta^2 \sum_{j=1}^n \|A_j B^\top B \nabla f(\theta_i^t)\|^2 & = \eta^2 (A\theta_i^t - \mathbf{b})^\top (AB^\top BA^\top)^\top AB^\top BA^\top (A\theta_i^t - \mathbf{b}) \leq 2\sigma_{\max}^2(Z)\eta^2 f(\theta_i^t), \\
\langle \nabla f(\theta_i^t) \mid B^\top B \nabla f(\theta_i^t) \rangle & \geq \sigma_{\min}(Z) f(\theta_i^t).
\end{aligned}$$

It remains to deal with the error terms separately. Observe that

$$\begin{aligned}
& \sum_{j=1}^n \eta \langle A_j \theta_i^t - b_j \mid A_j B^\top (E_i^t + \Delta G_i^t) \rangle \\
& \leq \eta \sum_{j=1}^n \left\{ \frac{\sigma_{\min}(Z)}{2} \|A_j \theta_i^t - b_j\|^2 + \frac{1}{2\sigma_{\min}(Z)} \|A_j B^\top (E_i^t + \Delta G_i^t)\|^2 \right\} \\
& = \sigma_{\min}(Z)\eta f(\theta_i^t) + \frac{\eta}{2\sigma_{\min}(Z)} \sum_{j=1}^n \|A_j B^\top (E_i^t + \Delta G_i^t)\|^2,
\end{aligned}$$

and

$$\begin{aligned}
& \sum_{j=1}^n \|A_j B^\top (E_i^t + \Delta G_i^t)\|^2 \\
& = \|AB^\top (E_i^t + \Delta G_i^t)\|^2 = (E_i^t + \Delta G_i^t)^\top BA^\top AB^\top (E_i^t + \Delta G_i^t) \\
& \leq \sigma_{\max}(Z) \|E_i^t + \Delta G_i^t\|^2 \leq 2\sigma_{\max}(Z) (\|E_i^t\|^2 + \|\Delta G_i^t\|^2)
\end{aligned}$$

Putting together into (31) and setting $\eta = \frac{\sigma_{\min}(Z)}{4\sigma_{\max}^2(Z)}$ yields

$$f(\theta_i^{t+1}) \leq \left(1 - \frac{\sigma_{\min}(Z)}{8\sigma_{\max}(Z)}\right) f(\theta_i^t) + \left(\eta^2 + \frac{\eta}{2\sigma_{\min}(Z)}\right) \cdot 2\sigma_{\max}(Z) (\|E_i^t\|^2 + \|\Delta G_i^t\|^2). \quad (32)$$

To bound $\|\Delta G_i^t\|^2$, we observe that by the algorithm and the initial condition $\theta_i^0 = \theta_j^0$, for any $t \geq 0$,

$$\begin{aligned}
\theta_j^{t+1} - \theta_i^{t+1} & = \theta_j^t - \theta_i^t - \eta(E_j^t - E_i^t) = -\eta \sum_{s=0}^t (E_j^s - E_i^s) \\
\|\theta_j^{t+1} - \theta_i^{t+1}\|^2 & \leq 2\eta^2 \sum_{s=0}^t [\|E_j^s\|^2 + \|E_i^s\|^2] \leq 4\eta^2 \sum_{s=0}^t \left(\frac{\bar{\epsilon}}{s+1}\right)^2 \leq \frac{2\eta^2 \pi^2 \bar{\epsilon}^2}{3}, \quad (33)
\end{aligned}$$

where the last inequality applied $\sum_{s=1}^{\infty} s^{-2} = \pi^2/6$, together with Fact 1 and the assumption $\epsilon_t = \bar{\epsilon}/(t+1)$. Then, the error ΔG_i^t can be bounded by

$$\begin{aligned}
\|\Delta G_i^t\|^2 & = \frac{1}{n^2} \left\| \sum_{j=1}^n BA_j^\top A_j (\theta_j^t - \theta_i^t) \right\|^2 \\
& \leq \frac{1}{n^2} \sum_{j=1}^n \|BA_j^\top A_j\|^2 \|\theta_j^t - \theta_i^t\|^2 \stackrel{(33)}{\leq} \frac{2\eta^2 \pi^2 \bar{\epsilon}^2}{3n^2} \sum_{j=1}^n \|BA_j^\top A_j\|^2. \quad (34)
\end{aligned}$$

With the notation $\sigma_{BA^\top A} := \frac{1}{n} \sum_{j=1}^n \|BA_j^\top A_j\|^2$, (32) gives us

$$\begin{aligned} f(\theta_i^{t+1}) &\stackrel{(34)}{\leq} \left(1 - \frac{1}{8\kappa^2}\right) f(\theta_i^t) + \left(\eta^2 + \frac{\eta}{2\sigma_{\min}(Z)}\right) \cdot 2\sigma_{\max}^2(Z) \left(\epsilon_t^2 + \frac{2\eta^2\pi^2\bar{\epsilon}^2}{3n}\sigma_{BA^\top A}\right) \\ &\leq \left(1 - \frac{1}{8\kappa^2}\right)^{t+1} f(\theta_i^0) + \sum_{s=0}^t \left(1 - \frac{\sigma_{\min}(Z)\eta}{2}\right)^{t-s} \left(\eta^2 + \frac{\eta}{2\sigma_{\min}(Z)}\right) \cdot 2\sigma_{\max}^2(Z) \left(\epsilon_s^2 + \frac{2\eta^2\pi^2\bar{\epsilon}^2}{3n}\sigma_{BA^\top A}\right) \\ &\leq \left(1 - \frac{1}{8\kappa^2}\right)^{t+1} f(\theta_i^0) + \frac{1}{6\sigma_{\min}^2(Z)} \left(1 + \frac{\sigma_{\min}^2(Z)}{4n\sigma_{\max}^4(Z)}\sigma_{BA^\top A}\right) (\sigma_{\min}^2(Z) + 2\sigma_{\max}^2(Z))\pi^2\bar{\epsilon}^2, \end{aligned}$$

where we have simplified notations by setting $\kappa = \frac{\sigma_{\max}(Z)}{\sigma_{\min}(Z)}$. By adjusting $\bar{\epsilon}$, we can achieve dimension-independent linear convergence.

Communication Complexity Let $\xi_t = \sum_{i=1}^n \|g_i^t - \frac{1}{n} \sum_{j=1}^n g_j^t\|^2$. We recall from Fact 1 that achieving a consensus error of $\bar{\epsilon}/(t+1)$ at iteration t via CHOCO-GOSSIP requires

$$\frac{82}{\delta^2\omega} \log \frac{(t+1)\xi_t}{\bar{\epsilon}} \text{ rounds of communication.}$$

Applying b -bits random quantization to our projected H -dimensional vector, we have

$$\frac{1}{\omega} = 1 + \min\{H/(2^b - 1)^2, \sqrt{H}/(2^b - 1)\}.$$

Fix $\epsilon > 0$, to find an ϵ -optimal solution, we set

$$\bar{\epsilon} := \sqrt{\epsilon \frac{3}{\pi^2(1+\kappa^2)} \left(1 + \frac{1}{4n\kappa^2\sigma_{\max}^2(Z)}\sigma_{BA^\top A}\right)^{-1}}$$

and $t \geq T_\epsilon := 8\kappa^2 \log(2f(\theta_i^0)/\epsilon)$. Altogether, the number of communication rounds for achieving an ϵ optimal solution is bounded by:

$$\begin{aligned} \frac{82}{\delta^2\omega} \sum_{t=1}^{T_\epsilon} \log \frac{(t+1)\xi_t}{\bar{\epsilon}} &\leq \frac{82T_\epsilon}{\delta^2\omega} \left(\max_{t \in [T_\epsilon]} \log(\xi_t) + \log(T_\epsilon + 1) + \log(1/\bar{\epsilon}) \right) \\ &= \mathcal{O} \left(\frac{\kappa^2}{\delta^2\omega} \left(\log \frac{1}{\epsilon} + \log \left(1 + \frac{1}{n\sigma_{\max}^2(Z)}\sigma_{BA^\top A}\right) \right) \log \frac{1}{\epsilon} \right) \end{aligned}$$

where we have assumed that $\max_{t \in [T_\epsilon]} \log(\xi_t)$ is dominated by $\max\{\log(1/\bar{\epsilon}), \log(T_\epsilon)\}$.

Each communication round requires to send $\mathcal{O}((b+1)H)$ bits per agent. We can optimize the choice of b by

$$\begin{aligned} b^* &= \arg \min_b bH\omega^{-1} = \arg \min_b \left\{ \min \{bH + bH^2/(2^b - 1)^2, bH + bH^{3/2}/(2^b - 1)\} \right\} \\ \Rightarrow b^*H\omega^{-1} &= \mathcal{O}(H \log H), \quad \text{with } b^* = \log(H^{1/2} + 1) \end{aligned}$$

Under the properties of A, B described in Proposition 1, by Lemma E.1, simplifying $\|BA_j^\top A_j\|$ shows that with probability $1 - \zeta$,

$$\|BA_j^\top A_j\| \leq \|A_j\|^2 \|B\| \leq \left(\sqrt{H} + \sqrt{D} + \sqrt{2 \log \frac{2}{\zeta}} \right)^3$$

As such, we have $\sigma_{BA^\top A} \leq \left(\sqrt{H} + \sqrt{D} + \sqrt{2 \log \frac{2}{\zeta}} \right)^6$. Putting together gives the communication complexity upper bound of \mathcal{C}_{op} and the proof is completed.

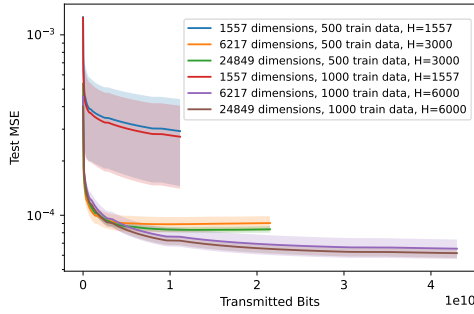


Figure 2: **Overparameterized Kernel Regression with Alg. 2.** Testing MSE against the number of bit transmitted on the whole network, averaged over 5 random seeded runs.

D	H	Compression Ratio
1557	1557	1.099
6217	3000	0.529
6217	6000	1.057
24849	3000	0.132
24849	6000	0.264

Table 1: **Communication Compression Ratio of Alg. 2** against vanilla DGD with double precision. 8-bits quantization is applied in CHOCO-GOSSIP and the average rounds of gossip T_g is 7.78.

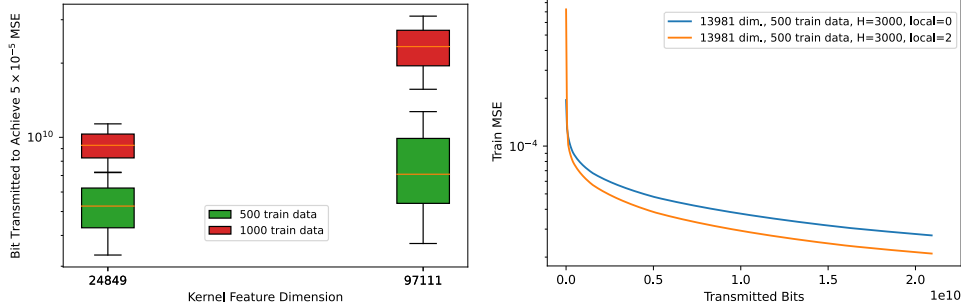


Figure 3: **Overparameterized Kernel Regression with Alg. 2.** (Left) Additional simulation results for large scale overparameterization ($19\times$). (Right) Additional simulation results of Alg. 2 combining with local updates, compared to when no local updates are applied.

F.3 Additional Numerical Result

Fig. 2 shows the test MSE against communication cost for the cases of $N = 2500$ and $N = 5000$. The models are evaluated on the testing dataset of 23175 samples. Our result indicates that increasing the dimension of kernel features decreases the testing MSE. It suggests that overparameterization improves generalization on unseen data. The communication budget limits the number of training samples, thus bottlenecks the generalization accuracy. In the scenario where both computation and communication budgets are limited, Fig. 2 shows that the generalization power of Alg. 2 benefits more from spending into communication budget (thus more training samples) than from overparametrizing, under the situation where $H \ll D$ and more training samples are available with no cost. Table 1 shows the compression ratio for the schemes used in our experiments.

Lastly, Fig. 3 (left) provides additional results for the large-scale overparameterized scenario, where $\frac{D}{N}$ is as large as 19 which is close to the degree of over-parameterization in practice. Comparing the two cases with $D = 24,849$, $D = 97,111$ where the latter case uses 4 times more parameters, we observe that the bit complexity to reach an MSE of 10^{-5} only increases by 1.4 times. Fig. 3 (right) examines empirically the effects of running $L_{\text{local}} = 2$ local steps on the train MSE performance against the number of transmitted bits.

G The Proof of Theorem 4.4

G.1 Additional Notations

Let us further define some notations before we go into the details of the proof. Denote the parameter in t -th iteration as $\theta^t = (W_l^t)_{l=1}^L$; $P_l = \text{vec}(O_l)$ as the vectorized output from layer l ; the pre-activation output from layer l as Q_l . Similarly, we can define $P_{l,k}, Q_{l,k}$ for each agent. Specifically, $O_L = P_L = Q_L$. Further, let us define some notations related to the singular values of the weight matrices.

$$\bar{\lambda}_l = \begin{cases} \frac{2}{3} (1 + s_{\max}(W_l^0)), & \text{for } l \in \{1, 2\} \\ s_{\max}(W_l^0), & \text{for } l \in \{3, \dots, L\} \end{cases}, \quad \underline{\lambda}_l = s_{\min}(W_l^0), \quad \underline{\lambda}_{i \rightarrow j} = \prod_{l=i}^j \underline{\lambda}_l, \quad \bar{\lambda}_{i \rightarrow j} = \prod_{l=i}^j \bar{\lambda}_l. \quad (35)$$

Specifically, $\bar{\lambda}_l$ represents the quantity related to the largest singular value of weight matrices for different layers; $\underline{\lambda}_l$ represents smallest singular value of weight matrices for different layers.

Further, denote $\lambda_O = s_{\min}(\tilde{B} \cdot a(XW_1^0))$ as the smallest singular value of the output from first hidden layer at initialization multiplied by \tilde{B} . Let us define some notations related to gradient. Recall we have defined $u_l = \text{vec}(\nabla_{W_l} f(\theta))$. Let us define the per iteration gradient as: $g_k^t = (Bu_{2,k}^t, u_{3,k}^t; \dots; u_{L,k}^t)$ that collects all the (compressed) gradient of each layer. Further, let us denote $\tilde{u}_{2,k}^t = B^\top Bu_{2,k}^t$, $\tilde{g}_k^t = (\tilde{u}_{2,k}^t; u_{3,k}^t; \dots; u_{L,k}^t)$, $q_k^t = (z_{2,k}^t; z_{3,k}^t \dots; z_{L,k}^t)$, $\tilde{q}_k^t = (B^\top z_{2,k}^t; z_{3,k}^t \dots; z_{L,k}^t)$, where $z_{2,k}^t$ quantizes $Bu_{2,k}^t$ and $z_{l,k}^t$ quantizes $u_{l,k}^t$ for $l \geq 3$; denote $\tilde{z}_{2,k}^t = B^\top z_{2,k}^t$. Note, that in the above notation, both quantities with $\tilde{\cdot}$ has the original dimension D . Further, \tilde{q}_k^t is the actual vector that gets used at k -th agent, while \tilde{g}_k^t represents a ‘‘virtual’’ vector which has not been quantized.

Now let us define:

$$\begin{aligned} u_l^t &:= \sum_{k=1}^K u_{l,k}^t, \quad u^t = (u_2^t, u_3^t, \dots, u_L^t), \\ g_k^t &:= \sum_{k=1}^K g_k^t = (Bu_{2,k}^t; u_{3,k}^t; \dots; u_{L,k}^t), \quad \tilde{g}^t := \sum_{k=1}^K \tilde{g}_k^t = (\tilde{u}_{2,k}^t; u_{3,k}^t; \dots; u_{L,k}^t), \\ q^t &:= \sum_{k=1}^K q_k^t = (z_2^t; z_3^t; \dots; z_L^t), \quad \tilde{q}^t := \sum_{k=1}^K \tilde{q}_k^t = (\tilde{z}_2^t; z_3^t; \dots, z_L^t). \end{aligned}$$

Further, we define $\Delta W_l^t, \tilde{\Delta} W_l^t$, which are the unvectorized \tilde{u}_2^t and \tilde{q}_2^t for $l = 2$ and unvectorized \tilde{u}_l^t and \tilde{q}_l^t for $l \geq 3$:

$$\text{vec}(\Delta W_l^t) = \begin{cases} \tilde{u}_2^t & l = 2 \\ u_l^t & l \geq 3, \end{cases} \quad \text{vec}(\tilde{\Delta} W_l^t) = \begin{cases} \tilde{z}_2^t & l = 2 \\ z_l^t & l \geq 3. \end{cases} \quad (36)$$

That is, $\Delta W_l^t, \tilde{\Delta} W_l^t \in \mathbb{R}^{n_{l-1} \times n_l}$.

Using the above notation, the ‘Update’ step in Algorithm 1, it can be expressed as

$$\theta^{t+1} = \theta^t - \eta \tilde{q}^t. \quad (37)$$

Denote $\Sigma_l^t = \text{diag}[\text{vec}(a'(Q_l^t))] \in \mathbb{R}^{N n_l \times N n_l}$, which is a diagonal matrix, whose diagonal entries are the vectorized gradient of the activation function in each layer. Recall that we have defined $B = \mathbf{I}_{n_2} \otimes \tilde{B}$ in Section 4.2, where \tilde{B} is a Gaussian random matrix of size $H \times n_1$. For convenience, let us assume $s_{\max}(B) = s_{\max}(\tilde{B}) \geq 1$.

G.2 Useful Lemma

Now we first state a collection of results from [Nguyen and Mondelli, 2020] that will be used in our proof. Notice that in the aforementioned work, the same pyramidal neural network structure and the l_2 loss function are used, so the loss function is the same as (6). It follows that all the properties of the loss function can be reused.

Lemma G.1. *Suppose Assumption 4 and 5 hold, for each θ^t , we have the following relations:*

$$1. u_l^t = \left(\mathbf{I}_{n_l} \otimes O_{l-1}^{t\top} \right) \prod_{p=l+1}^L \Sigma_{p-1}^t (W_p^t \otimes \mathbf{I}_N) (O_L^t - y), \quad (38)$$

$$2. \frac{\partial O_L^t}{\partial \text{vec}(W_l)} = \prod_{p=0}^{L-l-1} \left(W_{L-p}^{t\top} \otimes \mathbf{I}_N \right) \Sigma_{L-p-1}^t (\mathbf{I}_{n_l} \otimes F_{l-1}^t), \quad (39)$$

$$3. \|u_2^t\| \geq s_{\min} \left(O_1^{t\top} \right) \prod_{p=3}^L s_{\min} (\Sigma_{p-1}^t) s_{\min} (W_p^t) \|O_L^t - y\|, \quad (40)$$

$$4. \|u_l^t\| \leq \|X\|_F \prod_{\substack{p=1 \\ p \neq l}}^L s_{\max}(W_p^t) \|O_L^t - y\|, \quad (41)$$

$$5. \|g^t\| \leq s_{\max}(\tilde{B}) L \|X\|_F \frac{\prod_{l=1}^L s_{\max}(W_l^t)}{\min_{l \in [L]} s_{\max}(W_l^t)} \|O_L^t - y\|. \quad (42)$$

Furthermore, given with θ^a and θ^b , if $\bar{\Lambda}_l \geq \max(s_{\max}(W_l^a), s_{\max}(W_l^b))$ for some scalars $\bar{\Lambda}_l$. Let $\tilde{R} := \prod_{p=1}^L \max(1, \bar{\Lambda}_p)$. Then, for $l \in [L]$,

$$6. \|O_L^a - O_L^b\|_F \leq \sqrt{L} \|X\|_F \frac{\prod_{l=1}^L \bar{\Lambda}_l}{\min_{l \in [L]} \bar{\Lambda}_l} \|\theta^a - \theta^b\|, \quad (43)$$

$$7. \left\| \frac{\partial \text{vec}(O_L(\theta^a))}{\partial \text{vec}(W_l^a)} - \frac{\partial \text{vec}(O_L(\theta^b))}{\partial \text{vec}(W_l^b)} \right\|_2 \leq \sqrt{L} \|X\|_F \tilde{R} \left(1 + L\rho \|X\|_F \tilde{R} \right) \|\theta^a - \theta^b\|. \quad (44)$$

$$8. \left\| \text{vec}(\nabla f(\theta^a)) - \text{vec}(\nabla f(\theta^b)) \right\| \leq \left(L\sqrt{L} \|X\|_F^2 \frac{\prod_{l=1}^L \bar{\Lambda}_l^2}{\min_{l \in [L]} \bar{\Lambda}_l^2} + L\sqrt{L} \|X\|_F \tilde{R} (1 + L\rho \|X\|_F \tilde{R}) \right) \|\theta^a - \theta^b\|, \quad (45)$$

Now let us discuss the properties above one by one. The relations (38) and (39) show how the vectorized gradients of each layer, as well as the vectorized gradients of the output over each layer are computed, which are true regardless of the algorithm; see the first and the second equalities in [Nguyen and Mondelli, 2020, Lemma 4.1], respectively; (40) is the lower bound of the norm of the vectorized gradient over W_2 , which is true as long as the network has the pyramidal structure (and regardless of the algorithm); see the third relation in [Nguyen and Mondelli, 2020, Lemma 4.1]; (41) is the upper bound of the norm of vectorized gradient over W_l , which is true regardless of the algorithm; see the third relation in [Nguyen and Mondelli, 2020, Lemma 4.2]; (42) can be derived by summing over l in (41). To be specific, we have

$$\begin{aligned} \|g^t\| &\leq \|B u_2^t\| + \sum_{l=3}^L \|u_l^t\| \leq s_{\max}(\tilde{B}) \|u_2^t\| + \sum_{l=3}^L \|u_l^t\| \\ &\leq s_{\max}(\tilde{B}) \sum_{l=2}^L \|u_l^t\| \\ &\leq s_{\max}(\tilde{B}) \sum_{l=2}^L \|X\|_F \prod_{\substack{p=1 \\ p \neq l}}^L s_{\max}(W_p) \|O_L - y\| \end{aligned}$$

$$\leq s_{\max}(\tilde{B})L\|X\|_F \frac{\prod_{l=1}^L s_{\max}(W_l^t)}{\min_{l \in [L]} s_{\max}(W_l^t)} \|O_L^t - y\|.$$

The relation in (43) gives the upper bound of the gap between output layer with two sets of parameters θ^a and θ^b , and it is regardless of the algorithm; see [Nguyen and Mondelli, 2020, Eq. (19)]; Similarly, (44) states the gap between the vectorized Jacobian matrix over W_l ; see [Nguyen and Mondelli, 2020, Eq. (20)]; Finally, (45) computes the Lipschitz constants for the gradient, it is independent of the algorithm, and comes from plugging in (42), (43) and (44) into [Nguyen and Mondelli, 2020, Eq. (22)].

G.3 Initialization Strategy

As stated in Theorem 4.4, special initialization strategy is required. Now let us describe the initialization in detail. Recall the definition of $\bar{\lambda}_l$ and $\underline{\lambda}_l$ in (35) and $\lambda_O = s_{\min}(\tilde{B} \cdot a(XW_1^0))$. Initialize θ^0 such that the following holds:

$$\lambda_O^2 \geq \frac{48 \cdot 6^{L-2} \|X\|_F \bar{\lambda}_{1 \rightarrow L} s_{\max}^2(\tilde{B}) \|O_L^0 - y\|}{(\nu^2)^{L-2} \lambda_{3 \rightarrow L}^2} \max\left(\frac{1}{\lambda_2}, \max_{l \geq 3} \frac{2}{\lambda_l \lambda_l}\right). \quad (46)$$

To satisfy the above relation, it requires that $\|O_L^0 - y\|$ cannot be large, which means the θ^0 should not be far from the optimal solution so that the initial loss is small. Further, it requires the gap between $\bar{\lambda}_l$ and $\underline{\lambda}_l$ is small. From (45), we know that small $\bar{\lambda}_l$ induces small Lipschitz constant for gradient, so the condition requires that the Lipschitz constant is not very large, which further implies that the optimization landscape is smooth. On the other hand, from (40), we know the lower bound of the norm of u_2^t is related to $\underline{\lambda}_l$. So the initialization condition guarantees the lower bound of $\|u_2^t\|$ is not too small, which avoids vanished gradient.

The initialization can be realized by using the procedure suggested [Nguyen and Mondelli, 2020, Sec. 3.1]. The idea is that we can scale up W_1^0 to make λ_O not too small, and then randomly choose W_2^0 with small entries. Then for all $l \geq 3$, set W_l^0 as scaled identity matrices (top block as scaled identities) with large entries.

G.4 Formal Statement and Proof of Theorem 4.4

First, let us state the formal Theorem 4.4:

Theorem 4.4 *Consider using Alg. 1 to solve the problem (6), with X being full row rank. Suppose θ^0 is initialized as (46). Choose \tilde{B} such as $\text{rank}(\tilde{B}O_1^0) = N$; choose $\tilde{F}(\cdot)$ and $F(\cdot)$ as in (7). Set stepsize η and bits number b as following:*

$$\eta = \frac{\phi}{Q_0 \Lambda^2}, \quad b = \max\left(\log\left(\frac{1}{\tau} + 1\right), b_0\right),$$

where ϕ is defined in (47), τ is defined in (50), Λ is defined in (48), Q_0 is defined in (49). The following hold true:

$$f(\theta^{t+1}) \leq \left(1 - \frac{1}{4}\eta \cdot \phi\right) f(\theta^t).$$

To compute an ϵ -optimal solution (6), each agent is required to transmit:

$$\frac{4b}{\eta\phi} \cdot \left(Hn_2 + \sum_{l=2}^{L-1} n_l n_{l+1}\right) \log(f(\theta^0)/\epsilon) \quad \text{bits/agent}.$$

Proof. To begin with, let us set the following constants:

$$\alpha = 1 - \frac{1}{4}\eta\phi, \quad \gamma^t = \lambda\sqrt{(\alpha)^{t+1} f(\theta^0)}, \quad \eta = \frac{\phi}{Q_0 \Lambda^2}, \quad \lambda = 6\sqrt{2} s_{\max}(\tilde{B}) \eta Q_0 \Lambda, \quad (47)$$

$$\begin{aligned}
C_1 &:= \frac{\phi}{16\sqrt{2}Hn_1n_2^2\|\tilde{B}\|_\infty\lambda\|O_L^0 - y\|}, & C_2 &:= \frac{\phi \min_{l \geq 3} \lambda_l}{32\sqrt{2}\lambda n_2 n_3 \|O_L^0 - y\|}, \\
C_3 &:= \frac{1}{\sqrt{2}\eta H n_1 n_2^2 \|\tilde{B}\|_\infty \lambda \|O_L^0 - y\|}, & C_4 &:= \frac{\min_{l \geq 3} \lambda_l}{2\sqrt{2}\eta \lambda n_2 n_3 \|O_L^0 - y\|} \\
C_5 &:= \sqrt{\frac{\phi}{4\eta \left(\frac{1}{2} + \frac{1}{\beta}\right) \left(H n_1 n_2^2 \|\tilde{B}\|_\infty^2 + \sum_{l=3}^L n_{l-1}^2 n_l^2\right) \Lambda^2 \lambda^2}}, \\
C_6 &:= \frac{\sqrt{1 - \frac{1}{4}\eta\phi - \frac{1}{6}}}{1 + \eta s_{\max}(\tilde{B}) Q_0 \left(H n_2 + \sum_{l=3}^L n_{l-1} n_l\right)},
\end{aligned}$$

where

$$\phi = \left(\frac{1}{2}\nu\right)^{L-2} \lambda_{O\Delta_{3 \rightarrow L}}^2, \quad \Lambda = \left(\frac{3}{2}\right)^{L-1} \cdot L s_{\max}^2(\tilde{B}) \|X\|_F \frac{\bar{\lambda}_{1 \rightarrow L}}{\min_{l \geq 2} \bar{\lambda}_l} \quad (48)$$

$$Q_0 = L\sqrt{L} \left(\frac{3}{2}\right)^{2(L-1)} \|X\|_F^2 \frac{\prod_{l=1}^L \bar{\lambda}_l^2}{\min_{l \in [L]} \bar{\lambda}_l^2} + L\sqrt{L} \|X\|_F R (1 + L\rho \|X\|_F R), \quad (49)$$

$$R = s_{\max}^2(\tilde{B}) \prod_{p=1}^L \max\left(1, \frac{3}{2}\bar{\lambda}_p\right).$$

Define

$$\tau := \min(C_1, C_2, C_3, C_4, C_5, C_6). \quad (50)$$

Further, let us define

$$\Lambda_l := \left(\frac{3}{2}\right)^{L-1} \|X\|_F \frac{\bar{\lambda}_{1 \rightarrow L}}{\bar{\lambda}_l}. \quad (51)$$

For convenience, let us assume $Q_0 > 1$. Let us explain the above constants: α is the constants based on which the objective function contracts; η is the stepsize; ϕ is related to the lower bound of $\|u_2^t\|$; Λ is related to the upper bound of $\|\tilde{g}^t\|$; Q_0 is the Lipschitz constant for the full gradient; τ is the parameter related to quantization; Λ_l is related to the upper bound for $\|u_l^t\|$.

The majority of the proof consists of showing the following relations by induction:

$$\left\{ \begin{array}{l}
s_{\max}(W_l^t) \stackrel{(i)}{\leq} \frac{3}{2}\bar{\lambda}_l, l = \{2, 3, \dots, L\}, \\
s_{\min}(W_l^t) \stackrel{(ii)}{\geq} \frac{1}{2}\lambda_l, l \in \{3, \dots, L\}, \\
s_{\max}(W_2^t + \eta\tilde{\Delta}W_2^{t-1} - \eta\Delta W_2^{t-1}) \stackrel{(iii)}{\leq} \frac{3}{2}\bar{\lambda}_l, l = \{2, 3, \dots, L\}, \\
s_{\min}(W_2^t + \eta\tilde{\Delta}W_2^{t-1} - \eta\Delta W_2^{t-1}) \stackrel{(iv)}{\geq} \frac{1}{2}\lambda_l, l = \{3, \dots, L\}, \\
f(\theta^{t-1} - \eta\tilde{g}^{t-1}) \stackrel{(v)}{\leq} (1 - \eta\phi)f(\theta^{t-1}), \\
f(\theta^t) \stackrel{(vi)}{\leq} (\alpha)^t f(\theta^0), \\
\|Bu_2^t - z_2^t\|_\infty \stackrel{(vii)}{\leq} \tau\gamma^{t-1} = \tau\lambda\sqrt{(\alpha)^t f(\theta^0)}, \\
\|u_l^t - z_l^t\|_\infty \stackrel{(viii)}{\leq} \tau\gamma^{t-1} = \tau\lambda\sqrt{(\alpha)^t f(\theta^0)}.
\end{array} \right. \quad (52)$$

Let us explain the meanings of the above relations. Relations (i) and (ii) provide the upper and lower bounds of the singular values of weight matrices in each iteration; (iii) and (iv) give the upper and lower bound of the singular values of weight matrices after one step update without quantization; (v) shows the decrease of loss function after one step of update without quantization; (vi) shows the

linear decrease of loss function in each iteration; (vii) and (viii) provide the error bound of gradient after quantization.

Compared to the induction proof in [Nguyen and Mondelli, 2020], the key challenges are: (1) Our analysis includes the quantization of gradient; (2) The update of parameter is not the simple gradient but a function of gradient.

Our proof consists of two steps:

Step 1: We show the above relations hold for $t = 1$.

Step 2: We show that if the above relations hold for t , then they hold for $t + 1$.

We show Step 1 first. For simplicity, we set $g^0 = q^0$, it is easy to verify the relations hold if we choose q^0 in Algorithm 1 with large enough C . Step 1 will be shown in the following five substeps: (a) Show (i) and (ii) in (52); (b) Show (iii) and (iv) in (52); (c) Show (v) in (52); (d) Show (vi) in (52); (e) Show (vii) and (viii) in (52).

(Step 1.a) We will show

$$\begin{cases} s_{\max}(W_l^1) \leq \frac{3}{2}\bar{\lambda}_l, & l \in [L], l \in \{2, 3, \dots, L\} \\ s_{\min}(W_l^1) \geq \frac{1}{2}\underline{\lambda}_l, & l \in \{3, \dots, L\}. \end{cases}$$

We will use the fact that, according to the update rule in (37), we have $q^0 = g^0$.

For $l = 2$, we have:

$$\begin{aligned} \|W_2^1 - W_2^0\|_F &\stackrel{(i)}{=} \eta \|\tilde{z}_2^0\| = \eta \|B^\top B u_2^0\| \stackrel{(ii)}{\leq} \eta s_{\max}^2(\tilde{B}) \|u_2^0\| \\ &\stackrel{(iii)}{\leq} \eta s_{\max}^2(\tilde{B}) \Lambda_2 \|O_L^0 - y\| \\ &\stackrel{(iv)}{\leq} \sum_{t=0}^{\infty} \eta s_{\max}^2(\tilde{B}) \Lambda_2 \|O_L^0 - y\| \cdot \sqrt{\alpha}^t \\ &\leq \eta s_{\max}^2(\tilde{B}) \Lambda_2 \|O_L^0 - y\| \frac{1}{1 - \alpha} (1 + \alpha^{\frac{1}{2}}) \\ &\stackrel{(v)}{\leq} \frac{8s_{\max}^2(\tilde{B}) \Lambda_2}{\phi} \|O_L^0 - y\| \\ &\stackrel{(vi)}{\leq} \frac{1}{2}, \end{aligned}$$

where (i) is from the update rule in (37); (ii) extracts the largest singular value of B ; (iii) uses the upper bound of the gradient norm in (41) and the definition of Λ_2 in (51); (iv) uses the fact that $\sum_{t=0}^{\infty} \sqrt{\alpha} > 1$; (v) plugs in the definition of α and $1 + \alpha^{\frac{1}{2}} < 2$; (vi) is because the initialization strategy in (46).

Next, we show the case where $l \geq 3$. We have

$$\begin{aligned} \|W_l^1 - W_l^0\|_F &\stackrel{(i)}{=} \eta \|u_l^0\| \stackrel{(ii)}{\leq} \eta \Lambda_l \|u_l^0\| \\ &\stackrel{(iii)}{\leq} \sum_{t=0}^{\infty} \eta \Lambda_l \|O_L^0 - y\| \cdot \sqrt{\alpha}^t \\ &= \eta \Lambda_l \|O_L^0 - y\| \frac{1}{1 - \alpha^{\frac{1}{2}}} \\ &\leq \eta \Lambda_l \|O_L^0 - y\| \frac{1}{1 - \alpha} (1 + \alpha^{\frac{1}{2}}) \\ &\stackrel{(iv)}{\leq} \frac{8\Lambda_l}{\phi} \|O_L^0 - y\| \\ &\stackrel{(v)}{\leq} \frac{1}{4}\underline{\lambda}_l < \frac{1}{2}\underline{\lambda}_l, \end{aligned}$$

where (i) is because the update of parameter in (37); (ii) uses the upper bound of the gradient norm in (41) and definition of Λ_l in (51); (iii) uses the fact that $\sum_{t=0}^{\infty} \sqrt{\alpha} > 1$; (iv) plugs in the choice of α and $1 + \alpha^{\frac{1}{2}} < 2$; (v) comes from the initialization strategy in (46).

Applying Weyl' inequality to the matrices W_l^0 and $(W_l^1 - W_l^0)$, we have

$$\begin{cases} s_{\max}(W_l^1) \leq \bar{\lambda}_l + \frac{1}{2}\bar{\lambda}_l = \frac{3}{2}\bar{\lambda}_l, l \in \{3, \dots, L\}, \\ s_{\max}(W_l^1) \leq \bar{\lambda}_l + 1 = \frac{3}{2}\bar{\lambda}_l, l = 2, \\ s_{\min}(W_l^1) \geq \underline{\lambda}_l - \frac{1}{2}\bar{\lambda}_l = \frac{1}{2}\underline{\lambda}_l, l \in \{3, \dots, L\}. \end{cases}$$

This concludes the proof of this substep.

(Step 1.b) We will show

$$\begin{cases} s_{\max}(W_l^1 + \eta\tilde{\Delta}W_l^0 - \eta\Delta W_l^0) \leq \frac{3}{2}\bar{\lambda}_l, l = \{2, 3, \dots, L\}, \\ s_{\min}(W_l^1 + \eta\tilde{\Delta}W_l^0 - \eta\Delta W_l^0) \geq \frac{1}{2}\underline{\lambda}_l, l = \{3, \dots, L\}. \end{cases}$$

Since we have $\tilde{q}^0 = \tilde{g}^0$, it is easy to derive that

$$\|\eta\tilde{\Delta}W_l^0 - \eta\Delta W_l^0\|_F = \eta\|\tilde{\Delta}W_l^0 - \Delta W_l^0\| = \eta\|\tilde{q}^0 - \tilde{g}^0\| = 0.$$

Thus we can conclude

$$\begin{cases} s_{\max}(W_l^1 + \eta\tilde{\Delta}W_l^0 - \eta\Delta W_l^0) = s_{\max}(W_l^1) \leq \frac{3}{2}\bar{\lambda}_l, l = \{2, 3, \dots, L\} \\ s_{\min}(W_l^1 + \eta\tilde{\Delta}W_l^0 - \eta\Delta W_l^0) = s_{\min}(W_l^1) \geq \frac{1}{2}\underline{\lambda}_l, l = \{3, \dots, L\} \end{cases}$$

This concludes the proof for this substep.

(Step 1.c) We will show

$$f(\theta^0 - \eta\tilde{g}^0) \leq (1 - \eta\alpha_0)f(\theta^0).$$

From (Step 1.a) and (Step 1.b) we know

$$\max(s_{\max}(W_l^0), s_{\max}(W_l^1)) \leq \frac{3}{2}\bar{\lambda}_l.$$

Using the above relation, we can upper bound the differences of the gradients by using (45). More specifically, using the the definition of Q_0 in (49), we have

$$\|\text{vec}(\nabla f(\theta^1)) - \text{vec}(\nabla f(\theta^0))\| \leq Q_0\|\theta^1 - \theta^0\|, \quad (53)$$

where we repeat the definition of Q_0 below (where R is defined in (49)):

$$Q_0 = L\sqrt{L} \left(\frac{3}{2}\right)^{2(L-1)} \|X\|_F^2 \frac{\prod_{l=1}^L \bar{\lambda}_l^2}{\min_{l \in [L]} \bar{\lambda}_l^2} + L\sqrt{L}\|X\|_F R (1 + L\rho\|X\|_F R),$$

Further, it is easy to verify that for any $\hat{\theta}^0$ between θ^0 and θ^1 , we still have $s_{\max}(W_l(\hat{\theta}^0)) \leq \frac{3}{2}\bar{\lambda}_l$. So we can apply the same argument leading to (53), and obtain:

$$\|\text{vec}(\nabla f(\hat{\theta}^0)) - \text{vec}(\nabla f(\theta^0))\| \leq Q_0\|\hat{\theta}^0 - \theta^0\|, \forall \hat{\theta}^0 = \theta^0 + \delta(\theta^1 - \theta^0), \delta \in [0, 1]. \quad (54)$$

As long as for any $\hat{\theta}^0$ (54) holds, we can apply Decent Lemma. More specifically, we have

$$\begin{aligned} f(\theta^0 - \eta\tilde{g}^0) &\stackrel{(i)}{\leq} f(\theta^0) - \eta\langle u^0, \tilde{g}^0 \rangle + \frac{Q_0}{2}\eta^2\|\tilde{g}^0\|^2 \\ &\stackrel{(ii)}{=} f(\theta^0) - \eta\langle u_2^0, \tilde{u}_2^0 \rangle - \eta\sum_{l=3}^L \|u_l^0\|^2 + \frac{Q_0}{2}\eta^2\|\tilde{g}^0\|^2 \\ &\leq f(\theta^0) - \eta\langle u_2^0, \tilde{u}_2^0 \rangle + \frac{Q_0}{2}\eta^2\|\tilde{g}^0\|^2, \end{aligned} \quad (55)$$

where (i) uses Decent Lemma, (ii) uses the fact that u^0 is the stacked version of $\{u_i^0\}_{i \geq 2}$, and \tilde{g}^0 are stacked versions of $\{u_i^0\}_{i \geq 3}$ and \tilde{u}_2^0 .

Next, we will bound the inner product on the right hand side of the above relation:

$$\begin{aligned}
\langle u_2^0, \tilde{u}_2^0 \rangle &\stackrel{(i)}{=} \left\langle \left(\mathbf{I}_{n_2} \otimes O_1^{0\top} \right) \prod_{l=3}^L \Sigma_{l-1}^0 (W_l^0 \otimes \mathbf{I}_N) \left(O_L^{0\top} - y \right), B^\top B \left(\mathbf{I}_{n_2} \otimes O_1^{0\top} \right) \prod_{l=3}^L \Sigma_{l-1}^0 (W_l^0 \otimes \mathbf{I}_N) \left(O_L^{0\top} - y \right) \right\rangle \\
&\stackrel{(ii)}{=} \left\langle \left(\mathbf{I}_{n_2} \otimes O_1^{0\top} \right) \prod_{l=3}^L \Sigma_{l-1}^0 (W_l^0 \otimes \mathbf{I}_N) \left(O_L^{0\top} - y \right), \left(\mathbf{I}_{n_2} \otimes \tilde{B}^\top \tilde{B} O_1^{0\top} \right) \prod_{l=3}^L \Sigma_{l-1}^0 (W_l^0 \otimes \mathbf{I}_N) \left(O_L^{0\top} - y \right) \right\rangle \\
&\stackrel{(iii)}{=} \left\| \left(\mathbf{I}_{n_2} \otimes \tilde{B} O_1^{0\top} \right) \prod_{l=3}^L \Sigma_{l-1}^0 (W_l^0 \otimes \mathbf{I}_N) \left(O_L^0 - y \right) \right\|^2 \\
&\geq s_{\min}^2(\tilde{B} O_1^{0\top}) \prod_{l=3}^L s_{\min}^2(\Sigma_{l-1}^0) s_{\min}^2(W_l^0) \|O_L^0 - y\|^2 \\
&\stackrel{(iv)}{\geq} \left(\frac{1}{2} \nu \right)^{2(L-2)} \lambda_O^2 \lambda_{3 \rightarrow L}^2 \|O_L^0 - y\|^2, \tag{56}
\end{aligned}$$

where (i) uses the expression of u_2^0 by (38) and $\tilde{u}_2^0 = B^\top B u_2^0$; (ii) and (iii) are from the property of Kronecker product; (iv) is from Assumption 5 and fact that $s_{\max}(W_l^0) \geq \frac{1}{2} \lambda_l$.

Next, we upper bound the term $\|\tilde{g}^0\|^2$ in (55). We have

$$\begin{aligned}
\|\tilde{g}^0\|^2 &\stackrel{(i)}{=} \|B^\top B u_2^0\|^2 + \sum_{l=3}^L \|u_l^0\|^2 \stackrel{(ii)}{\leq} s_{\max}^4(\tilde{B}) \|u_2^0\|^2 + \sum_{l=3}^L \|u_l^0\|^2 \\
&\stackrel{(iii)}{\leq} s_{\max}^2(\tilde{B}) (\|B u_2^0\|^2 + \sum_{l=3}^L \|u_l^0\|^2) = s_{\max}^2(\tilde{B}) \|g^0\|^2 \\
&\stackrel{(iv)}{\leq} s_{\max}^2(\tilde{B}) \left(L s_{\max}(\tilde{B}) \|X\|_F \frac{\prod_{l=1}^L s_{\max}(W_l^0)}{\min_l s_{\max}(W_l^0)} \right)^2 \|O_L^0 - y\|^2 \\
&\stackrel{(v)}{\leq} \left(\frac{3}{2} \right)^{2(L-1)} L^2 s_{\max}^4(\tilde{B}) \|X\|_F^2 \frac{\bar{\lambda}_{1 \rightarrow L}^2}{\min_l \lambda_l^2} \|O_L^0 - y\|^2, \\
&\stackrel{(vi)}{=} \Lambda^2 \|O_L^0 - y\|^2 \tag{57}
\end{aligned}$$

where (i) uses the definition of \tilde{g}^0 ; (ii) extracts the largest singular value of B ; (iii) is because we assume $s_{\max}(\tilde{B}) \geq 1$; (iv) uses the upper bound of gradient in (42), (v) comes from the fact that $s_{\max}(W_l^0) \geq \frac{1}{2} \lambda_l$; (iii) is because the definition of Λ in (48).

Recall that we have defined

$$\phi := \left(\frac{1}{2} \nu \right)^{2(L-2)} \lambda_O^2 \lambda_{3 \rightarrow L}^2.$$

If we choose $\eta = \frac{\phi}{Q_0 \Lambda^2}$, then plug (56) and (57) into (55), we obtain

$$f(\theta^0 - \eta \tilde{g}^0) = f(\theta^0) - 2\eta \phi f(\theta^0) + \eta \phi f(\theta^0) = (1 - \eta \phi) f(\theta^0). \tag{58}$$

(Step 1.d) We will show

$$f(\theta^1) \leq \alpha f(\theta^0).$$

Since $\tilde{g}^0 = \tilde{q}^0$, we have $f(\theta^1) = f(\theta^0 - \eta \tilde{g}^0) \stackrel{(i)}{\leq} (1 - \eta \phi) f(\theta^0) < \alpha f(\theta^0)$, where (i) is because (58). This completes the proof of this step.

(Step 1.e) We will show that:

$$\begin{cases} \|B u_l^1 - v_l^1\|_\infty \leq \tau \gamma^0, & l = 2 \\ \|u_l^1 - v_l^1\|_\infty \leq \tau \gamma^0, & l = \{3, \dots, L\}. \end{cases}$$

Notice that for $l = 2$, we have

$$\begin{aligned}
\|Bu_2^1 - z_2^0\|_\infty &\stackrel{(i)}{\leq} \|Bu_2^1 - Bu_2^0\| + \|Bu_2^0 - z_2^0\|_\infty \stackrel{(ii)}{\leq} s_{\max}(\tilde{B})\|u_2^1 - u_2^0\| \\
&\stackrel{(iii)}{\leq} s_{\max}(\tilde{B})\|\text{vec}(\nabla f(\theta^1)) - \text{vec}(\nabla f(\theta^0))\| \\
&\stackrel{(iv)}{\leq} s_{\max}(\tilde{B})Q_0\|\theta^1 - \theta^0\| \stackrel{(v)}{=} \eta s_{\max}(\tilde{B})Q_0\|\tilde{q}^0\| \\
&\stackrel{(vi)}{=} \eta s_{\max}(\tilde{B})Q_0\|\tilde{g}^0\| \stackrel{(vii)}{\leq} \eta s_{\max}(\tilde{B})Q_0\Lambda\|O_L^0 - y\| \\
&\stackrel{(viii)}{\leq} \lambda\sqrt{\alpha f(\theta^0)} = \gamma^0,
\end{aligned}$$

where (i) uses triangle inequality; (ii) extracts the largest singular value of B and uses $g^0 = q^0$; (iii) is because $\|u_2^1 - u_2^0\| \leq \|u^1 - u^0\|$; (iv) is from (53), which uses (45) and definition of Q_0 in (49); (v) uses the update rule in (37); (vi) is because $g^0 = q^0$; (vii) uses the definition of Λ in (48); (viii) uses the definition of λ in (47) and the fact $\alpha = 1 - \frac{\phi^2}{4Q_0\Lambda^2} > 1 - \frac{1}{4} = \frac{3}{4}$.

Then by Lemma 4.1, with the correspondence that $c = Bu_2^1, p = z_2^0, r = \gamma^0$, we can obtain

$$\|z_2^1 - Bu_2^1\| \stackrel{(i)}{=} \|\text{quant}(Bu_2^1, z_2^0, \gamma^0, b) - Bu_2^1\|_\infty \stackrel{(ii)}{\leq} \tau\gamma^0,$$

where (i) comes from the ‘Quantize’ step in Algorithm 1; (ii) is from the first relation in Lemma 4.1.

For $l \geq 3$, we have

$$\begin{aligned}
\|u_l^1 - v_l^0\|_\infty &\stackrel{(i)}{\leq} \|u_l^1 - u_l^0\| + \|u_l^0 - z_l^0\|_\infty \stackrel{(ii)}{\leq} \|\text{vec}(\nabla f(\theta^1)) - \text{vec}(\nabla f(\theta^0))\| \\
&\stackrel{(iii)}{\leq} Q_0\|\theta^1 - \theta^0\| \stackrel{(iv)}{=} \eta Q_0\|\tilde{q}^0\| \stackrel{(v)}{=} \eta Q_0\|\tilde{g}^0\| \stackrel{(vi)}{\leq} \eta Q_0\Lambda\|O_L^0 - y\| \\
&\stackrel{(vii)}{\leq} \lambda\sqrt{\alpha f(\theta^0)} = \gamma_0,
\end{aligned}$$

where (i) uses triangle inequality; (ii) is because $g^0 = q^0$ and uses the fact $\|u_l^0 - z_l^0\| \leq \|g^0 - q^0\|$; (iii) is from (53); (iv) uses the update rule in (37); (v) is because $\tilde{g}^0 = \tilde{q}^0$; (vi) uses the upper bound of $\|\tilde{g}^0\|$ in (42) and the definition of Λ in (48); (vii) comes from the choice of λ in (47) and the fact $\alpha = 1 - \frac{\phi^2}{4Q_0\Lambda^2} > 1 - \frac{1}{4} = \frac{3}{4}$.

Then by Lemma 4.1, with the correspondence that $c = u_l^1, p = z_l^0, r = \gamma^0$, we can obtain

$$\|z_l^{t+1} - u_l^{t+1}\| \stackrel{(i)}{=} \|\text{quant}(u_l^1, z_l^0, \gamma^0, b) - u_l^1\|_\infty \stackrel{(ii)}{\leq} \tau\gamma^0.$$

where (i) comes from the ‘Quantize’ step in Algorithm 1; (ii) is from the first relation in Lemma 4.1. Now we have showed all the induction assumptions hold for $t = 1$.

Next, we show Step 2. Given these inequalities in (52) hold for iteration t , we aim to show that they hold for $t + 1$. We prove Step 2 in five substeps similarly as in the proof of Step 1.

(Step 2.a) We will show that

$$\begin{cases} s_{\max}(W_l^{t+1}) \leq \frac{3}{2}\bar{\lambda}_l, & l \in [L], l \in \{2, 3, \dots, L\} \\ s_{\min}(W_l^{t+1}) \geq \frac{1}{2}\underline{\lambda}_l, & l \in \{3, \dots, L\}. \end{cases}$$

For $l = 2$, we have:

$$\begin{aligned}
\|W_2^{t+1} - W_2^0\|_F &\stackrel{(i)}{\leq} \sum_{r=0}^t \|W_2^{r+1} - W_2^r\|_F \stackrel{(ii)}{\leq} \sum_{r=0}^t \|\tilde{z}_2^r\| \\
&\stackrel{(iii)}{\leq} \eta \sum_{r=0}^t (\|\tilde{u}_2^r\| + \|\tilde{u}_2^r - \tilde{z}_2^r\|) \\
&\stackrel{(iv)}{\leq} \eta \sum_{r=0}^t (s_{\max}^2(\tilde{B})\Lambda_2\|O_L^r - y\| + n_1 n_2 \|B^\top(Bu_2^r - z_2^r)\|_\infty)
\end{aligned}$$

$$\begin{aligned}
&\stackrel{(v)}{\leq} \eta \sum_{r=0}^t (s_{\max}^2(\tilde{B})\Lambda_2 \|O_L^r - y\| + Hn_1n_2^2 \|\tilde{B}^\top\|_\infty \|Bu_2^r - z_2^r\|_\infty) \\
&\stackrel{(vi)}{\leq} \eta \sum_{r=0}^t (s_{\max}^2(\tilde{B})\Lambda_2 \|O_L^r - y\| + Hn_1n_2^2 \|\tilde{B}^\top\|_\infty \tau \lambda \sqrt{(\alpha)^r f(\theta^0)}) \\
&\stackrel{(vii)}{=} (\eta s_{\max}^2(\tilde{B})\Lambda_2 + \frac{\sqrt{2}}{2} \eta Hn_1n_2^2 \|\tilde{B}^\top\|_\infty \tau \lambda) \|O_L^0 - y\| \sum_{r=0}^t \alpha^{\frac{r}{2}} \\
&= (\eta s_{\max}^2(\tilde{B})\Lambda_2 + \frac{\sqrt{2}}{2} \eta Hn_1n_2^2 \|\tilde{B}^\top\|_\infty \tau \lambda) \|O_L^0 - y\| \frac{1 - \alpha^{\frac{t+1}{2}}}{1 - \alpha^{\frac{1}{2}}} \\
&\stackrel{(viii)}{\leq} (\eta s_{\max}^2(\tilde{B})\Lambda_2 + \frac{\sqrt{2}}{2} \eta Hn_1n_2^2 \|\tilde{B}^\top\|_\infty \tau \lambda) \|O_L^0 - y\| \frac{1}{1 - \alpha} (1 + \alpha^{\frac{1}{2}}) \\
&\stackrel{(ix)}{\leq} \frac{8(s_{\max}^2(\tilde{B})\Lambda_2 + \frac{\sqrt{2}}{2} Hn_1n_2^2 \|\tilde{B}^\top\|_\infty \tau \lambda)}{\phi} \|O_L^0 - y\| \\
&\stackrel{(x)}{\leq} \frac{1}{4} + \frac{1}{4} = \frac{1}{2} < 1, \tag{59}
\end{aligned}$$

where (i) uses triangle inequality; (ii) uses the update rule in (37); (iii) uses triangle inequality; (iv) extracts the largest singular value of $B^\top B$, the upper bound of $\|u_2^l\|$ in (41), definition of Λ_2 in (51) and the relationship between l_2 and l_∞ norm; (v) uses the fact that $\|B^\top (Bu_2^r - z_2^r)\|_\infty \leq Hn_2 \|\tilde{B}\|_\infty \|Bu_2^r - z_2^r\|_\infty$; (vi) uses the induction assumption $\|Bu_2^r - z_2^r\|_\infty \leq \tau \gamma^{r-1}$ and definition of γ^{r-1} ; (vii) reorganizes the terms; (viii) is because $0 < \alpha < 1$; (ix) uses $1 + \alpha^{\frac{1}{2}} < 2$ and plugs in the choice of α ; (x) is because the initialization strategy in (46) and choice of τ in (50).

So Applying Weyl' inequality to the matrices W_2^t and $(W_2^{t+1} - W_2^t)$ and (59), we have

$$s_{\max}(W_2^{t+1}) \leq \bar{\lambda}_2 + 1 = \frac{3}{2} \bar{\lambda}_2.$$

For $l \geq 3$, we have the following relation:

$$\begin{aligned}
\|W_l^{t+1} - W_l^0\|_F &\stackrel{(i)}{\leq} \sum_{r=0}^t \|W_l^{r+1} - W_l^r\|_F \stackrel{(ii)}{\leq} \sum_{r=0}^t \|z_l^r\| \\
&\stackrel{(iii)}{\leq} \eta \sum_{r=0}^t (\|u_l^r\| + \|u_l^r - z_l^r\|) \\
&\stackrel{(iv)}{\leq} \eta \sum_{r=0}^t (\Lambda_l \|O_L^r - y\| + n_{l-1}n_l \|u_l^r - z_l^r\|_\infty) \\
&\stackrel{(v)}{\leq} \eta \sum_{r=0}^t (\Lambda_l \|O_L^r - y\| + n_{l-1}n_l \tau \lambda \sqrt{(\alpha)^r f(\theta^0)}) \\
&\stackrel{(vi)}{\leq} (\eta \Lambda_l + \frac{\sqrt{2}}{2} \eta \tau \lambda n_{l-1}n_l) \|O_L^0 - y\| \sum_{r=0}^t \alpha^{\frac{r}{2}} \\
&\stackrel{(vii)}{\leq} (\eta \Lambda_l + \frac{\sqrt{2}}{2} \eta \tau \lambda n_{l-1}n_l) \|O_L^0 - y\| \frac{1}{1 - \alpha} \\
&= (\eta \Lambda_l + \frac{\sqrt{2}}{2} \eta \tau \lambda n_{l-1}n_l) \|O_L^0 - y\| \frac{1}{1 - \alpha} (1 + \alpha^{\frac{1}{2}}) \\
&\stackrel{(viii)}{\leq} \frac{8(\Lambda_l + \frac{\sqrt{2}}{2} \tau \lambda n_{l-1}n_l)}{\phi} \|O_L^0 - y\| \\
&\stackrel{(ix)}{\leq} \frac{1}{8} \lambda_l + \frac{1}{8} \lambda_l = \frac{1}{4} \lambda_l < \frac{1}{2} \bar{\lambda}_l, \tag{60}
\end{aligned}$$

where (i) uses triangle inequality; (ii) uses the update rule in (37); (iii) uses triangle inequality; (iv) uses the upper bound of $\|u_l^l\|$ in (41), the definition of Λ_l in (51) and the relationship between l_2 and

l_∞ norm; (v) uses the induction assumption $\|u_l^r - z_l^r\| \leq \tau\gamma^{r-1}$; (vi) uses reorganizes the terms; (vii) is because $0 < \alpha < 1$; (viii) uses the fact $1 + \alpha^{\frac{1}{2}} < 2$ and plugs in the choice of α ; (ix) is from the initialization strategy in (46) and the choice of τ in (50).

Applying Weyl' inequality to the matrices W_l^{t+1} and $(W_l^{t+1} - W_l^t)$, we have

$$s_{\max}(W_l^{t+1}) \leq \bar{\lambda}_l + \frac{1}{4}\bar{\lambda}_l = \frac{5}{4}\bar{\lambda}_l < \frac{3}{2}\bar{\lambda}_l, \quad s_{\min}(W_l^{t+1}) \geq \bar{\lambda}_l - \frac{1}{2}\lambda_l = \frac{1}{2}\bar{\lambda}_l. \quad (61)$$

(Step 2.b) We will show that

$$\begin{cases} s_{\max}(W_l^{t+1} + \eta\tilde{\Delta}W_l^t - \eta\Delta W_l^t) \leq \frac{3}{2}\bar{\lambda}_l, l = \{2, 3, \dots, L\} \\ s_{\min}(W_l^{t+1} + \eta\tilde{\Delta}W_l^t - \eta\Delta W_l^t) \geq \frac{1}{2}\lambda_l, l = \{3, \dots, L\}. \end{cases}$$

For $l = 2$, we have

$$\begin{aligned} s_{\max}(W_2^{t+1} + \eta\tilde{\Delta}W_2^t - \eta\Delta W_2^t) &\stackrel{(i)}{\leq} s_{\max}(W_2^{t+1}) + \eta\|\tilde{\Delta}W_2^t - \Delta W_2^t\|_F \\ &\stackrel{(ii)}{=} s_{\max}(W_2^{t+1}) + \eta\|\tilde{z}_2^t - \tilde{u}_2^t\| \\ &\stackrel{(iii)}{\leq} s_{\max}(W_2^{t+1}) + \eta n_1 n_2 \|B^\top (Bu_2^r - z_2^r)\|_\infty \\ &\stackrel{(iv)}{\leq} s_{\max}(W_2^{t+1}) + H n_1 n_2^2 \|\tilde{B}^\top\|_\infty \|Bu_2^r - z_2^r\|_\infty \\ &\stackrel{(v)}{\leq} s_{\max}(W_2^{t+1}) + H n_1 n_2^2 \|\tilde{B}^\top\|_\infty \eta \tau \lambda \sqrt{(\alpha)^r f(\theta^0)} \\ &\stackrel{(vi)}{\leq} s_{\max}(W_2^{t+1}) + H n_1 n_2^2 \|\tilde{B}^\top\|_\infty \eta \tau \lambda \sqrt{f(\theta^0)} \stackrel{(v)}{\leq} \bar{\lambda}_l + \frac{1}{2} + \frac{1}{2} = \frac{3}{2}\bar{\lambda}_2, \end{aligned}$$

where (i) uses Weyl's inequality and the relationship between the Frobenius norm and l_2 norm of matrix; (ii) uses the definition in (36); (iii) uses the relationship between l_2 and l_∞ norm and the definition of \tilde{u}_2^t and \tilde{z}_2^t ; (iv) uses the fact that $\|B^\top (Bu_2^r - z_2^r)\|_\infty \leq H n_2 \|\tilde{B}\|_\infty \|Bu_2^r - z_2^r\|_\infty$ and the definition of γ^{r-1} ; (v) is from the induction assumption $\|Bu_2^r - z_2^r\|_\infty \leq \tau\gamma^{r-1}$; (vi) is from the fact that $\alpha \leq 1$; (v) is from the (59) and choice of τ and λ in (50) and (47).

For $l \geq 3$, by Weyl's inequality, we have

$$\begin{aligned} s_{\max}(W_l^{t+1} + \eta\tilde{\Delta}W_l^t - \eta\Delta W_l^t) &\stackrel{(i)}{\leq} s_{\max}(W_l^{t+1}) + \eta\|\tilde{\Delta}W_l^t - \Delta W_l^t\|_F \\ &\stackrel{(ii)}{=} s_{\max}(W_l^{t+1}) + \eta\|z_l^t - u_l^t\| \stackrel{(iii)}{\leq} s_{\max}(W_l^{t+1}) + \eta n_{l-1} n_l \tau \lambda \sqrt{(\alpha)^t f(\theta^0)} \\ &\stackrel{(iv)}{\leq} s_{\max}(W_l^{t+1}) + \eta n_{l-1} n_l \tau \lambda \sqrt{f(\theta^0)} \stackrel{(v)}{\leq} \bar{\lambda}_l + \frac{1}{4}\lambda_l + \frac{1}{4}\lambda_l \leq \frac{3}{2}\bar{\lambda}_l, \end{aligned} \quad (62)$$

where (i) uses the Weyl's inequality on W_l^{t+1} and $\eta(\tilde{\Delta}W_l^t - \Delta W_l^t)$; (ii) uses the update rule in (37); (iii) uses the relationship between l_2 and l_∞ norm and the induction assumption; (iv) is because $\alpha < 1$; (v) uses (60) and choice of τ and λ in (50) and (47).

Similarly,

$$\begin{aligned} s_{\min}(W_l^{t+1} + \eta\tilde{\Delta}W_l^t - \eta\Delta W_l^t) &\stackrel{(i)}{\geq} s_{\min}(W_l^{t+1}) - \eta\|\tilde{\Delta}W_l^t - \Delta W_l^t\|_F \\ &\stackrel{(ii)}{=} s_{\min}(W_l^{t+1}) - \eta\|z_l^t - u_l^t\| \stackrel{(iii)}{\geq} s_{\min}(W_l^{t+1}) - \eta n_{l-1} n_l \tau \lambda \sqrt{(\alpha)^t f(\theta^0)} \\ &\stackrel{(iv)}{\geq} s_{\min}(W_l^{t+1}) - \eta n_{l-1} n_l \tau \lambda \sqrt{f(\theta^0)} \stackrel{(v)}{\geq} \frac{3}{4}\lambda_l - \frac{1}{4}\lambda_l = \frac{1}{2}\lambda_l, \end{aligned}$$

where (i) uses Weyl's inequality on W_l^{t+1} and $\eta(\tilde{\Delta}W_l^t - \Delta W_l^t)$; (ii) plugs in the update rule in (37); (iii) uses the relationship between l_2 and l_∞ norm and the induction assumption; (iv) is because $\alpha < 1$; (v) comes from (60) and choice of τ and λ in (50) and (47).

(Step 2.c) We will show that:

$$f(\theta^t - \eta\tilde{g}^t) \leq (1 - \eta\phi)f(\theta^t).$$

From (1) and (2) we know

$$\max (s_{\max}(W_l^t), s_{\min}(W_l^t)) \leq \frac{3}{2}\bar{\lambda}_l.$$

Using the above relation, we can upper bound the differences of the gradients by using (45). More specifically, using the definition of Q_0 in (49), we have

$$\|\text{vec}(\nabla f(\theta^{t+1})) - \text{vec}(\nabla f(\theta^t))\| \leq Q_0 \|\theta^{t+1} - \theta^t\|,$$

where Q_0 is defined in (49).

Further, it is easy to verify that for any $\hat{\theta}^t$ between θ^t and θ^{t+1} , we still have $s_{\max}(W_l(\hat{\theta}^t)) \leq \frac{3}{2}\bar{\lambda}_l$. So we can apply the same argument leading to (53), and obtain:

$$\|\text{vec}(\nabla f(\hat{\theta}^t)) - \text{vec}(\nabla f(\theta^t))\| \leq Q_0 \|\hat{\theta}^t - \theta^t\|, \forall \hat{\theta}^t = \theta^t + \delta(\theta^{t+1} - \theta^t), \delta \in [0, 1].$$

So we have

$$\begin{aligned} f(\theta^t - \eta\tilde{g}^t) &\stackrel{(i)}{\leq} f(\theta^t) - \eta\langle u^t, \tilde{g}^t \rangle + \frac{Q_0}{2}\eta^2\|\tilde{g}^t\|^2 \\ &\stackrel{(ii)}{=} f(\theta^t) - \eta\langle u_2^t, \tilde{u}_2^t \rangle - \eta \sum_{l=3}^L \|u_l^t\|^2 + \frac{Q_0}{2}\eta^2\|\tilde{g}^t\|^2 \\ &\leq f(\theta^t) - \eta\langle u_2^t, \tilde{u}_2^t \rangle + \frac{Q_0}{2}\eta^2\|\tilde{g}^t\|^2, \end{aligned}$$

where (i) uses Decent Lemma; (ii) expands $\langle u^t, \tilde{g}^t \rangle$ by utilizing the stacked structure. Similar to (56), with the induction assumption $s_{\min}(W_l^t) \geq \frac{1}{2}\lambda_l, l \geq 3$, we have the following relation:

$$\langle u_2^t, \tilde{u}_2^t \rangle \geq \phi \|O_L^t - y\|^2.$$

Similar to (57), with the induction assumption $s_{\max}(W_l^t) \leq \frac{3}{2}\bar{\lambda}_l, l \geq 2$, we have

$$\|\tilde{g}^t\|^2 \leq \Lambda^2 \|O_L^t - y\|^2.$$

If we choose $\eta \leq \frac{\phi}{Q_0\Lambda^2}$, we obtain

$$f(\theta^t - \eta\tilde{g}^t) \leq f(\theta^t) - 2\eta\phi f(\theta^t) + \eta\phi f(\theta^t) = (1 - \eta\phi)f(\theta^t).$$

(Step 2.d) We will show that:

$$f(\theta^{t+1}) \leq \alpha f(\theta^t).$$

We have the following relation:

$$\begin{aligned} f(\theta^{t+1}) &= f(\theta^{t+1}) - f(\theta^t - \eta\tilde{g}^t) + f(\theta^t - \eta\tilde{g}^t) \\ &\stackrel{(i)}{\leq} \frac{1}{2}\|G(\theta^{t+1}) - G(\theta^t - \eta\tilde{g}^t)\|^2 + \langle G(\theta^{t+1}) - G(\theta^t - \eta\tilde{g}^t), G(\theta^t - \eta\tilde{g}^t) - y \rangle \\ &\quad + f(\theta^t + \eta\tilde{q}^t - \eta\tilde{g}^t) \\ &\stackrel{(ii)}{\leq} \frac{1}{2}\|G(\theta^{t+1}) - G(\theta^t - \eta\tilde{g}^t)\|^2 + \frac{1}{\beta}\|G(\theta^{t+1}) - G(\theta^t - \eta\tilde{g}^t)\|^2 + \frac{\beta}{2}\|G(\theta^t - \eta\tilde{g}^t) - y\|^2 \\ &\quad + f(\theta^t - \eta\tilde{g}^t) \\ &\stackrel{(iii)}{\leq} \left(\frac{1}{2} + \frac{1}{\beta}\right)\|G(\theta^{t+1}) - G(\theta^t - \eta\tilde{g}^t)\|^2 + (1 + \beta)(1 - \eta\phi)f(\theta^t) \\ &\stackrel{(iv)}{=} \left(\frac{1}{2} + \frac{1}{\beta}\right)\langle \text{vec}(\nabla G(\hat{\theta}^t)), -\eta\tilde{q}^t + \eta\tilde{g}^t \rangle^2 + (1 + \beta)(1 - \eta\phi)f(\theta^t) \\ &\stackrel{(v)}{\leq} \eta^2\left(\frac{1}{2} + \frac{1}{\beta}\right)\|\text{vec}(\nabla G(\hat{\theta}^t))\|^2\|\tilde{q}^t - \tilde{g}^t\|^2 + (1 + \beta)(1 - \eta\phi)f(\theta^t), \end{aligned}$$

where (i) expands and reorganizes the loss function; (ii) uses Young's inequality with constant β ; (iii) uses the induction assumption $f(\theta^{t+1}) \leq (1 - \eta\phi)f(\theta^t)$; (iv) uses the mean value Theorem where $\hat{\theta}^t$ will be discussed in the next paragraph, the update rule in (37); (v) uses Cauchy-Schwartz inequality.

Now we discuss $\hat{\theta}^t$ and derive a bound for $\|\text{vec}(\nabla G(\hat{\theta}^t))\|$. By the Mean Value Theorem we have

$$\hat{\theta}^t = \theta^{t+1} + \delta(\eta\tilde{g}^t - \eta\tilde{q}^t),$$

for some $\delta \in [0, 1]$. By (61) and (62), we know that

$$s_{\max}(W_l(\theta^{t+1})) \leq \frac{5}{4}\bar{\lambda}_l, \quad s_{\max}(W_l(\theta^t - \eta\tilde{g}^t)) \leq \frac{3}{2}\bar{\lambda}_l.$$

So it is easy to conclude that for $l \geq 3$,

$$\begin{aligned} s_{\max}(W_l(\hat{\theta}^t)) &\stackrel{(i)}{\leq} s_{\max}(W_l(\theta^{t+1})) + \delta\eta\|z_l^t - u_l^t\| \\ &\stackrel{(ii)}{\leq} s_{\max}(W_l(\theta^{t+1})) + \eta\|z_l^t - u_l^t\| \stackrel{(iii)}{\leq} \frac{3}{2}\bar{\lambda}_l, \end{aligned}$$

where (i) uses Wely's inequality on W_l^{t+1} and $\delta\eta(\tilde{\Delta}W_l^t - \Delta W_l^t)$; (ii) is because $\delta \in [0, 1]$; (iii) uses (62). We can derive the similar result for $l = 2$. By (42) without $\|O_L^t - y\|$, it is clear that

$$\|\text{vec}(\nabla G(\hat{\theta}^t))\| \leq \Lambda.$$

So we have

$$\begin{aligned} &\eta^2\left(\frac{1}{2} + \frac{1}{\beta}\right)\|\text{vec}(\nabla G(\hat{\theta}^t))\|^2\|\tilde{q}^t - \tilde{g}^t\|^2 + (1 + \beta)(1 - \eta\phi)f(\theta^t) \\ &\stackrel{(i)}{\leq} \eta^2\left(\frac{1}{2} + \frac{1}{\beta}\right)\Lambda^2(\|\tilde{u}_2^t - \tilde{v}_2^t\|^2 + \sum_{l=3}^L\|u_l^t - z_l^t\|^2) + (1 + \beta)(1 - \eta\phi)f(\theta^t) \\ &\stackrel{(ii)}{\leq} \eta^2\left(\frac{1}{2} + \frac{1}{\beta}\right)\Lambda^2(n_1^2n_2^2\|\tilde{u}_2^t - \tilde{z}_2^t\|_\infty^2 + \sum_{l=3}^L n_{l-1}^2n_l^2\|u_l^t - z_l^t\|_\infty^2) + (1 + \beta)(1 - \eta\phi)f(\theta^t) \\ &\stackrel{(iii)}{\leq} \eta^2\left(\frac{1}{2} + \frac{1}{\beta}\right)\Lambda^2(H^2n_1^4n_2^2\|\tilde{B}^\top\|_\infty^2\|Bu_2^t - z_2^t\|_\infty^2 + \sum_{l=3}^L n_{l-1}^2n_l^2\|u_l^t - z_l^t\|_\infty^2) + (1 + \beta)(1 - \eta\phi)f(\theta^t) \\ &\stackrel{(iv)}{\leq} \eta^2\left(\frac{1}{2} + \frac{1}{\beta}\right)\Lambda^2\left(H^2n_1^4n_2^2\|\tilde{B}^\top\|_\infty^2\tau^2\lambda^2(\alpha)^t f(\theta^0) + \sum_{l=3}^L n_{l-1}^2n_l^2\tau^2\lambda^2(\alpha)^t f(\theta^0)\right) + (1 + \beta)(1 - \eta\phi)f(\theta^t) \\ &= \left((1 + \beta)(1 - \eta\phi) + \eta^2\left(\frac{1}{2} + \frac{1}{\beta}\right)(H^2n_1^4n_2^2\|\tilde{B}^\top\|_\infty^2 + \sum_{l=3}^L n_{l-1}^2n_l^2)\Lambda^2\tau^2\lambda^2\right)(\alpha)^t f(\theta^0), \end{aligned}$$

(i) uses (42) with the fact that $s_{\max}(W_l(\hat{\theta}^t)) \leq \frac{3}{2}\bar{\lambda}_l$, and expands $\|\tilde{q}^t - \tilde{g}^t\|^2$; (ii) uses the relationship between l_2 and l_∞ norm; (iii) uses the fact that $\|B^\top(Bu_2^t - z_2^t)\|_\infty \leq Hn_2\|\tilde{B}\|_\infty\|Bu_2^t - z_2^t\|_\infty$; (iv) uses the induction assumption $\|Bu_2^t - z_2^t\| \leq \tau\gamma^{t-1}$, $f(\theta^t) \leq (\alpha)^t f(\theta^0)$ and the definition of γ^{t-1} .

Let $\beta = \frac{\frac{1}{2}\eta\phi}{1-\eta\phi}$, we have

$$\begin{aligned} &\left((1 + \beta)(1 - \eta\phi) + \eta^2\left(\frac{1}{2} + \frac{1}{\beta}\right)(Hn_1n_2^2\|\tilde{B}^\top\|_\infty^2 + \sum_{l=3}^L n_{l-1}^2n_l^2)\Lambda^2\tau^2\lambda^2\right)(\alpha)^t f(\theta^0) \\ &\stackrel{(i)}{\leq} \left(1 - \frac{1}{2}\eta\phi + \eta^2\left(\frac{1}{2} + \frac{2(1 - \eta\phi)}{\eta\phi}\right)(Hn_1n_2^2\|\tilde{B}^\top\|_\infty^2 + \sum_{l=3}^L n_{l-1}^2n_l^2)\Lambda^2\tau^2\lambda^2\right)(\alpha)^t f(\theta^0) \\ &\stackrel{(ii)}{\leq} \left(1 - \frac{1}{4}\eta\phi\right)(\alpha)^t f(\theta^0) = (\alpha)^{t+1}f(\theta^0), \end{aligned}$$

where (i) plugs in the choice of β ; (ii) is because the choice of τ and λ in (50) and (47).

(Step 2.e) We will show that:

$$\begin{cases} \|Bu_l^{t+1} - z_l^{t+1}\|_\infty \leq \tau\gamma^t, \quad l = 2 \\ \|u_l^{t+1} - z_l^{t+1}\|_\infty \leq \tau\gamma^t, \quad l = \{3, \dots, L\}. \end{cases}$$

For $l = 2$, we have

$$\begin{aligned}
& \|Bu_2^{t+1} - z_2^t\|_\infty \stackrel{(i)}{\leq} \|Bu_2^{t+1} - Bu_2^t\| + \|Bu_2^t - z_2^t\|_\infty \stackrel{(ii)}{\leq} s_{\max}(\tilde{B})\|u_2^{t+1} - u_2^t\| + \|Bu_2^t - z_2^t\|_\infty \\
& \stackrel{(iii)}{\leq} s_{\max}(\tilde{B})\|u^{t+1} - u^t\| + \|Bu_2^t - z_2^t\|_\infty \\
& \stackrel{(iv)}{\leq} s_{\max}(\tilde{B})\|\text{vec}(\nabla f(\theta^{t+1})) - \text{vec}(\nabla f(\theta^t))\| + \|Bu_2^t - z_2^t\|_\infty \\
& \stackrel{(v)}{\leq} s_{\max}(\tilde{B})Q_0\|\theta^{t+1} - \theta^t\| + \|Bu_2^t - z_2^t\|_\infty \\
& \stackrel{(vi)}{=} \eta s_{\max}(\tilde{B})Q_0\|\tilde{q}^t\| + \|Bu_2^t - z_2^t\|_\infty \\
& \stackrel{(vii)}{\leq} \eta s_{\max}(\tilde{B})Q_0(\|\tilde{g}^t\| + \|\tilde{g}^t - \tilde{q}^t\|) + \|Bu_2^t - z_2^t\|_\infty \\
& \stackrel{(viii)}{\leq} \eta s_{\max}(\tilde{B})Q_0\Lambda\|O_L^t - y\| + \eta s_{\max}^2(\tilde{B})Q_0\|g^t - q^t\| + \|Bu_2^t - z_2^t\|_\infty \\
& \stackrel{(ix)}{\leq} \eta s_{\max}(\tilde{B})Q_0\Lambda\|O_L^t - y\| + \eta s_{\max}^2(\tilde{B})Q_0 \sum_{l=3}^L \|u_l^t - z_l^t\| + \eta s_{\max}^2(\tilde{B})Q_0\|Bu_2^t - z_2^t\| + \|Bu_2^t - z_2^t\|_\infty \\
& \stackrel{(x)}{\leq} \eta s_{\max}(\tilde{B})Q_0\Lambda\|O_L^t - y\| + \eta s_{\max}^2(\tilde{B})Q_0 H n_2 \|u_2^t - z_2^t\|_\infty + \eta s_{\max}^2(\tilde{B})Q_0 \sum_{l=3}^L n_{l-1} n_l \|u_l^t - z_l^t\|_\infty + \|Bu_2^t - z_2^t\|_\infty \\
& \stackrel{(xi)}{\leq} \left(1 + \eta s_{\max}^2(\tilde{B})Q_0(H n_2 + \sum_{l=3}^L n_{l-1} n_l)\right) \tau \gamma^{t-1} + \eta s_{\max}(\tilde{B})Q_0\Lambda\|O_L^t - y\| \\
& \stackrel{(xii)}{\leq} \lambda \sqrt{(\alpha)^t f(\theta^0)} \left(\tau + \tau \eta s_{\max}^2(\tilde{B})Q_0(H n_2 + \sum_{l=3}^L n_{l-1} n_l) + \frac{\sqrt{2} \eta s_{\max}(\tilde{B})Q_0\Lambda}{\lambda} \right) \\
& \stackrel{(xiii)}{\leq} \lambda \sqrt{\alpha} \sqrt{(\alpha)^t f(\theta^0)} = \lambda \gamma^t,
\end{aligned}$$

where (i) uses triangle inequality; (ii) extracts the largest singular value of B , (iii) uses the fact $\|u^{t+1} - u^t\| \geq \|u_2^{t+1} - u_2^t\|$; (iv) uses the definition of u^t ; (v) upper bounds the differences of the gradients by using (45) (vi) uses the update rule in (37); (vii) uses triangle inequality; (viii) uses the upper bound of gradient norm in (42), and it extracts the largest singular value of B and definition of Λ in (48); (ix) uses the stacked structure of g^t and q^t to expand $\|g^t - q^t\|$; (x) uses the relationship between l_2 and l_∞ norm; (xi) reorganizes the terms and uses induction assumption; (xii) uses the definition of γ^{t-1} ; (xiii) is from the choice of τ and λ in (50) and (47).

Then by Lemma 4.1, with the correspondence that $c = Bu_2^{t+1}$, $p = z_2^t$, $r = \gamma^t$, we can obtain

$$\|z_2^{t+1} - Bu_2^{t+1}\| \stackrel{(i)}{=} \|\text{quant}(Bu_2^{t+1}, z_2^t, \gamma^t, b) - Bu_2^{t+1}\|_\infty \stackrel{(ii)}{\leq} \tau \gamma^t.$$

where (i) comes from the ‘Quantize’ step in Algorithm 1; (ii) is from the first relation in Lemma 4.1.

For $l \geq 3$, we have

$$\begin{aligned}
& \|u_l^{t+1} - z_l^t\|_\infty \stackrel{(i)}{\leq} \|u_l^{t+1} - u_l^t\| + \|u_l^t - z_l^t\|_\infty \stackrel{(ii)}{\leq} \|u^{t+1} - u^t\| + \|u_l^t - z_l^t\|_\infty \\
& \stackrel{(iii)}{=} \|\text{vec}(\nabla f(\theta^{t+1})) - \text{vec}(\nabla f(\theta^t))\| + \|u_l^t - z_l^t\|_\infty \\
& \stackrel{(iv)}{\leq} \eta Q_0\|\tilde{q}^t\| + \|u_l^t - z_l^t\|_\infty \\
& \stackrel{(v)}{\leq} \eta Q_0(\|\tilde{g}^t\| + \|\tilde{g}^t - \tilde{q}^t\|) + \|u_l^t - z_l^t\|_\infty \\
& \stackrel{(vi)}{\leq} \eta Q_0\Lambda\|O_L^t - y\| + \eta s_{\max}(\tilde{B})Q_0\|g^t - q^t\| + \|u_l^t - z_l^t\|_\infty \\
& \stackrel{(vii)}{\leq} \eta Q_0\Lambda\|O_L^t - y\| + \eta s_{\max}(\tilde{B})Q_0 \sum_{l=3}^L \|u_l^t - z_l^t\| + \eta s_{\max}(\tilde{B})Q_0\|Bu_2^t - z_2^t\| + \|u_l^t - z_l^t\|_\infty
\end{aligned}$$

$$\begin{aligned}
&\stackrel{(viii)}{\leq} \eta Q_0 \Lambda \|O_L^t - y\| + \eta s_{\max}(\tilde{B}) Q_0 H n_2 \|B u_2^t - z_2^t\|_\infty + \eta s_{\max}(\tilde{B}) Q_0 \sum_{l=3}^L n_{l-1} n_l \|u_l^t - z_l^t\|_\infty + \|u_2^t - z_2^t\|_\infty \\
&\stackrel{(ix)}{\leq} \left(1 + \eta s_{\max}(\tilde{B}) Q_0 (H n_2 + \sum_{l=3}^L n_{l-1} n_l)\right) \tau \gamma^{t-1} + \eta Q_0 \Lambda \|O_L^t - y\| \\
&\stackrel{(x)}{\leq} \lambda \sqrt{f(\theta^t)} \left(\tau + \tau \eta s_{\max}(\tilde{B}) Q_0 (H n_2 + \sum_{l=3}^L n_{l-1} n_l) + \frac{\sqrt{2} \eta Q_0 \Lambda}{\lambda} \right) \\
&\stackrel{(xi)}{\leq} \lambda \sqrt{\alpha} \sqrt{f(\theta^t)} = \lambda \gamma^t,
\end{aligned}$$

where (i) uses triangle inequality; (ii) uses the fact $\|u^{t+1} - u^t\| \geq \|u_2^{t+1} - u_2^t\|$; (iii) uses the definition of u^t ; (iv) uses the update rule in (37); (v) uses triangle inequality; (vi) uses the upper bound of gradient norm in (42), extracts the largest singular value of B and uses the definition of Λ in (48); (vii) uses the stacked structure of g^t and q^t to expand $\|g^t - q^t\|$; (viii) uses the relationship between l_2 and l_∞ norm; (ix) uses the induction assumption; (x) uses the definition of γ^{t-1} ; (xi) is because the choice of τ and λ in (50) and (47).

Then by Lemma 4.1, with the correspondence that $c = B u_2^{t+1}$, $p = z_2^t$, $r = \gamma^t$, we can obtain

$$\|z_i^{t+1} - u_i^{t+1}\| \stackrel{(i)}{=} \|\text{quant}(u_i^{t+1}, z_i^t, \gamma^t, b) - B u_i^{t+1}\|_\infty \stackrel{(ii)}{\leq} \tau \gamma^t,$$

where (i) comes from the ‘Quantize’ step in Algorithm 1; (ii) is from the first relation in Lemma 4.1.

Now we have proved that (52) holds. So for $t > 0$, there is

$$f(\theta^t) \leq (\alpha)^t f(\theta^0), \quad \text{where} \quad \alpha = 1 - \frac{1}{4} \eta \phi = 1 - \frac{\phi^2}{4 Q_0 \Lambda^2}.$$

Thus, if we want the objective function to compute an ϵ -optimal solution, the total number of iterations is $\log(f(\theta^0)/\epsilon)/\log(1/(1 - \frac{\phi^2}{4 Q_0 \Lambda^2}))$. Since in each iteration, each agent k transmits a length- H vector q_k^t , so we conclude that the total number of bits each node needs to communicate is $\log(f(\theta^0)/\epsilon)/\log(1/(1 - \frac{\phi^2}{4 Q_0 \Lambda^2}))$ bits. Notice that $\log(1/(1 - \frac{\phi^2}{4 Q_0 \Lambda^2})) = -\log(1 - \frac{\phi^2}{4 Q_0 \Lambda^2}) \sim \frac{4 Q_0 \Lambda^2}{\phi^2}$, so we can derive the simplified total number of bits as $b \cdot \frac{4 Q_0 \Lambda^2}{\phi^2} \cdot \log(f(\theta^0)/\epsilon)$. \square