

# Basis Matters: Better Communication-Efficient Second Order Methods for Federated Learning

Xun Qian<sup>1,2</sup>   Rustem Islamov<sup>1,3</sup>   Mher Safaryan<sup>1</sup>   Peter Richtárik<sup>1</sup>

<sup>1</sup>KAUST   <sup>2</sup>JD Explore Academy   <sup>3</sup>Institut Polytechnique de Paris

## The Problem

$$\min_{x \in \mathbb{R}^d} f(x) \stackrel{\text{def}}{=} \frac{1}{n} \sum_{i=1}^n f_i(x), \quad (1)$$

where each function  $f_i : \mathbb{R}^d \rightarrow \mathbb{R}$  represents the local loss associated with the data owned by device or client  $i \in [n] \stackrel{\text{def}}{=} \{1, 2, \dots, n\}$  only.

## Algorithm

- $[\cdot]_\mu$ : the projection on the set  $\{\mathbf{A} \in \mathbb{R}^{d \times d} \mid \mathbf{A} = \mathbf{A}^\top, \mathbf{A} \succeq \mu \mathbf{I}\}$ .

**Algorithm 1:** Basis Learn with [Bidirectional Compression](#) (BL1)

**Parameters:** Hessian learning rate  $\alpha \geq 0$ ; [model learning rate](#)  $\eta \geq 0$ ;

[gradient compression probability](#)  $p \in (0, 1]$ ; compression operators

$\{\mathcal{C}_1^k, \dots, \mathcal{C}_n^k\}$  and  $\mathcal{Q}^k$ ; Basis  $\{\mathbf{B}_i^{jl}\}$  in  $\mathbb{R}^{d \times d}$  for each  $i$

**Initialization:**  $x^0 = w^0 = z^0 \in \mathbb{R}^d$ ;  $\mathbf{L}_i^0 \in \mathbb{R}^{d \times d}$ ,  $\mathbf{H}_i^0 = \sum_{jl} (\mathbf{L}_i^0)_{jl} \mathbf{B}_i^{jl}$ , and

$\mathbf{H}^0 \stackrel{\text{def}}{=} \frac{1}{n} \sum_{i=1}^n \mathbf{H}_i^0$ ;  $\xi^0 = 1$

**for each device**  $i = 1, \dots, n$  **in parallel do**

  if  $\xi^k = 1$

$w^{k+1} = z^k$ , compute  $\nabla f_i(z^k)$  and send to the server

  if  $\xi^k = 0$

$w^{k+1} = w^k$

  Compute  $\nabla^2 f_i(z^k)$  and send  $\mathbf{S}_i^k \stackrel{\text{def}}{=} \mathcal{C}_i^k(h^i(\nabla^2 f_i(z^k)) - \mathbf{L}_i^k)$  to the server.

  Update local Hessian shifts  $\mathbf{L}_i^{k+1} = \mathbf{L}_i^k + \alpha \mathbf{S}_i^k$ ,

$\mathbf{H}_i^{k+1} = \mathbf{H}_i^k + \alpha \sum_{jl} (\mathbf{S}_i^k)_{jl} \mathbf{B}_i^{jl}$

**end**

**on server**

  if  $\xi^k = 1$

$w^{k+1} = z^k$ ,  $g^k = \nabla f(z^k)$

  if  $\xi^k = 0$

$w^{k+1} = w^k$ ,  $g^k = [\mathbf{H}^k]_\mu (z^k - w^k) + \nabla f(w^k)$

$x^{k+1} = z^k - [\mathbf{H}^k]_\mu^{-1} g^k$      $\mathbf{H}^{k+1} = \mathbf{H}^k + \frac{\alpha}{n} \sum_{i=1}^n \sum_{jl} (\mathbf{S}_i^k)_{jl} \mathbf{B}_i^{jl}$

  Send  $v^k \stackrel{\text{def}}{=} \mathcal{Q}^k(x^{k+1} - z^k)$  to all devices  $i \in [n]$

  Update the model  $z^{k+1} = z^k + \eta v^k$

  Send  $\xi^{k+1} \sim \text{Bernoulli}(p)$  to all devices  $i \in [n]$

**for each device**  $i = 1, \dots, n$  **in parallel do**

  | Update the model  $z^{k+1} = z^k + \eta v^k$

**end**

## Compressor

- $\mathcal{C} : \mathbb{R}^{d \times d} \rightarrow \mathbb{R}^{d \times d}$  is called a *contraction compressor* if there exists a constant  $0 < \delta \leq 1$  such that

$$\mathbb{E} [\|\mathbf{A} - \mathcal{C}(\mathbf{A})\|_F^2] \leq (1 - \delta) \|\mathbf{A}\|_F^2, \quad \forall \mathbf{A} \in \mathbb{R}^{d \times d}. \quad (2)$$

- $\mathcal{C} : \mathbb{R}^{d \times d} \rightarrow \mathbb{R}^{d \times d}$  is an *unbiased compressor* if there exists a constant  $\omega \geq 0$  such that for any  $\mathbf{A} \in \mathbb{R}^{d \times d}$

$$\mathbb{E} [\mathcal{C}(\mathbf{A})] = \mathbf{A} \text{ and } \mathbb{E} [\|\mathcal{C}(\mathbf{A})\|_F^2] \leq (\omega + 1) \|\mathbf{A}\|_F^2. \quad (3)$$

## Assumptions

**Assumption 1** (i)  $\mathcal{Q}^k$  ( $\mathcal{Q}_i^k$ ) is an unbiased compressor with parameter  $\omega_M$  and  $0 < \eta \leq 1/(\omega_M + 1)$ . (ii) For all  $j \in [d]$ ,  $(z^k)_j$  in Algorithm 1 (  $(z^k)_j$  in Algorithm 2 ) is a convex combination of  $\{(x^t)_j\}_{t=0}^k$  for  $k \geq 0$ .

**Assumption 2** (i)  $\mathcal{Q}^k$  ( $\mathcal{Q}_i^k$ ) is a contraction compressor with parameter  $\delta_M$  and  $\eta = 1$ . (ii)  $\mathcal{Q}^k$  ( $\mathcal{Q}_i^k$ ) is deterministic, i.e.,  $\mathbb{E}[\mathcal{Q}^k(x)] = \mathcal{Q}^k(x)$  for any  $x \in \mathbb{R}^d$ .

**Assumption 3** (i)  $\mathcal{C}_i^k$  is an unbiased compressor with parameter  $\omega$  and  $0 < \alpha \leq 1/(\omega + 1)$ .

(ii) For all  $i \in [n]$  and  $j, l \in [d]$ ,  $(\mathbf{L}_i^k)_{jl}$  is a convex combination of  $\{h^i(\nabla^2 f_i(z^t))_{jl}\}_{t=0}^k$  in Algorithm 1 (  $\{h^i(\nabla^2 f_i(z^t))_{jl}\}_{t=0}^k$  in Algorithm 2 ) for  $k \geq 0$ .

**Assumption 4** (i)  $\mathcal{C}_i^k$  is a contraction compressor with parameter  $\delta$  and  $\alpha = 1$ . (ii)  $\mathcal{C}_i^k$  is deterministic, i.e.,  $\mathbb{E}[\mathcal{C}_i^k(\mathbf{A})] = \mathcal{C}_i^k(\mathbf{A})$  for any  $\mathbf{A} \in \mathbb{R}^{d \times d}$ .

**Assumption 5** We have  $\|\nabla^2 f_i(x) - \nabla^2 f_i(y)\| \leq H\|x - y\|$ ,  $\|\nabla^2 f_i(x) - \nabla^2 f_i(y)\|_F \leq H_1\|x - y\|$ ,  $\|h^i(\nabla^2 f_i(x)) - h^i(\nabla^2 f_i(y))\|_F \leq M_1\|x - y\|$ ,  $\max_{jl} \{|h^i(\nabla^2 f_i(x))_{jl} - h^i(\nabla^2 f_i(y))_{jl}|\} \leq M_2\|x - y\|$ ,  $\max_{jl} \{\|\mathbf{B}_i^{jl}\|_F\} \leq R$  for any  $x, y \in \mathbb{R}^d$  and  $i \in [n]$ . For Algorithm 2, we assume each  $f_i$  is  $\mu$ -strongly convex.

## Some Notations

$$N_B \stackrel{\text{def}}{=} \begin{cases} 1 & \text{if the bases } \{\mathbf{B}_i^{jl}\}_{j,l \in [d]} \text{ are all orthogonal} \\ d^2 & \text{otherwise} \end{cases} \quad (4)$$

$$(A_M, B_M) \stackrel{\text{def}}{=} \begin{cases} (\eta, \eta) & \text{if Asm. 1(ii) holds} \\ \left(\frac{\delta_M}{4}, \frac{6}{\delta_M} - \frac{7}{2}\right) & \text{if Asm. 2(ii) holds} \end{cases} \quad (5)$$

$$(A, B) \stackrel{\text{def}}{=} \begin{cases} (\alpha, \alpha) & \text{if Asm. 3(i) holds} \\ \left(\frac{\delta}{4}, \frac{6}{\delta} - \frac{7}{2}\right) & \text{if Asm. 4(i) holds} \end{cases} \quad (6)$$

For any  $k \geq 0$ , denote  $\mathcal{H}^k \stackrel{\text{def}}{=} \frac{1}{n} \sum_{i=1}^n \|\mathbf{L}_i^k - \mathbf{L}^*\|_F^2$ ,  $\Phi_1^k \stackrel{\text{def}}{=} \|z^k - x^*\|^2 + \frac{A_M(1-p)}{2p} \|w^k - x^*\|^2$ , where  $\mathbf{L}_i^* \stackrel{\text{def}}{=} h^i(\nabla^2 f_i(x^*))$ ,  $\Phi_2^k \stackrel{\text{def}}{=} \mathcal{H}^k + \frac{4BM_M^2}{A_M} \|x^k - x^*\|^2$ .

## Convergence Result

**Linear convergence of BL1:** Let Assumption 5 hold. Let Assumption 1 (i) or Assumption 2 (i) hold. Assume  $\|z^k - x^*\|^2 \leq \frac{A_M \mu^2}{4H^2 B_M}$  and  $\mathcal{H}^k \leq \frac{A_M \mu^2}{16N_B R^2 B_M}$  for  $k \geq 0$ . Thenxfor  $k \geq 0$  we have

$$\mathbb{E}[\Phi_1^k] \leq \left(1 - \frac{\min\{A_M, p\}}{2}\right)^k \Phi_1^0.$$

**Superlinear convergence of BL1:** Let  $\eta = 1$ ,  $\xi^k \equiv 1$  and  $\mathcal{Q}^k(x) \equiv x$  for any  $x \in \mathbb{R}^d$  and  $k \geq 0$ . Let Assumption 5 hold. Let Assumption 3 (i) or Assumption 4 (i) hold. Assume  $\|z^k - x^*\|^2 \leq \frac{A_M \mu^2}{4H^2 B_M}$  and  $\mathcal{H}^k \leq \frac{A_M \mu^2}{16N_B R^2 B_M}$  for  $k \geq 0$ . Then we have

$$\mathbb{E}[\Phi_2^k] \leq \theta_1^k \Phi_2^0,$$

$$\mathbb{E} \left[ \frac{\|x^{k+1} - x^*\|^2}{\|x^k - x^*\|^2} \right] \leq \theta_1^k \left( \frac{A_M H^2}{8B_M \mu^2} + \frac{2N_B R^2}{\mu^2} \right) \Phi_2^0,$$

for  $k \geq 0$ , where  $\theta_1 \stackrel{\text{def}}{=} \left(1 - \frac{\min\{4A, A_M\}}{4}\right)$ .

## Algorithm

- $[\mathbf{A}]_s = (\mathbf{A} + \mathbf{A}^\top)/2$  for any  $\mathbf{A} \in \mathbb{R}^{d \times d}$ .

**Algorithm 2:** Basis Learn with [Bidirectional Compression and Partial Participation](#) (BL2)

**Parameters:**  $\alpha > 0$ ;  $\eta > 0$ ; matrix compression operators  $\{\mathcal{C}_1^k, \dots, \mathcal{C}_n^k\}$ ;

$p \in (0, 1]$ ;  $0 < \tau \leq n$

**Initialization:**  $w_i^0 = z_i^0 = x^0 \in \mathbb{R}^d$ ;  $\mathbf{L}_i^0 \in \mathbb{R}^{d \times d}$ ;  $\mathbf{H}_i^0 = \sum_{jl} (\mathbf{L}_i^0)_{jl} \mathbf{B}_i^{jl}$ ;

$l_i^0 = \|[\mathbf{H}_i^0]_s - \nabla^2 f_i(w_i^0)\|_F$ ;  $g_i^0 = ([\mathbf{H}_i^0]_s + l_i^0 \mathbf{I}) w_i^0 - \nabla f_i(w_i^0)$ ; Moreover:

$\mathbf{H}^0 = \frac{1}{n} \sum_{i=1}^n \mathbf{H}_i^0$ ;  $l^0 = \frac{1}{n} \sum_{i=1}^n l_i^0$ ;  $g^0 = \frac{1}{n} \sum_{i=1}^n g_i^0$

**on server**

$x^{k+1} = ([\mathbf{H}^k]_s + l^k \mathbf{I})^{-1} g^k$ , choose a subset  $S^k \subseteq [n]$  such that

$\mathbb{P}[i \in S^k] = \tau/n$  for all  $i \in [n]$

$v_i^k = \mathcal{Q}_i^k(x^{k+1} - z_i^k)$ ,  $z_i^{k+1} = z_i^k + \eta v_i^k$  for  $i \in S^k$

$z_i^{k+1} = z_i^k$ ,  $w_i^{k+1} = w_i^k$  for  $i \notin S^k$

  Send  $v_i^k$  to the selected devices  $i \in S^k$

**for each device**  $i = 1, \dots, n$  **in parallel do**

**for participating devices**  $i \in S^k$  **do**

$z_i^{k+1} = z_i^k + \eta v_i^k$ ,  $\mathbf{S}_i^k \stackrel{\text{def}}{=} \mathcal{C}_i^k(h^i(\nabla^2 f_i(z_i^{k+1})) - \mathbf{L}_i^k)$

$\mathbf{L}_i^{k+1} = \mathbf{L}_i^k + \alpha \mathbf{S}_i^k$ ,  $\mathbf{H}_i^{k+1} = \mathbf{H}_i^k + \alpha \sum_{jl} (\mathbf{S}_i^k)_{jl} \mathbf{B}_i^{jl}$

$l_i^{k+1} = \|[\mathbf{H}_i^{k+1}]_s - \nabla^2 f_i(z_i^{k+1})\|_F$

    Sample  $\xi_i^{k+1} \sim \text{Bernoulli}(p)$

**if**  $\xi_i^k = 1$

$w_i^{k+1} = z_i^{k+1}$ ,  $g_i^{k+1} = ([\mathbf{H}_i^{k+1}]_s + l_i^{k+1} \mathbf{I}) w_i^{k+1} - \nabla f_i(w_i^{k+1})$ , send

$g_i^{k+1} - g_i^k$  to server

**if**  $\xi_i^k = 0$

$w_i^{k+1} = w_i^k$ ,  $g_i^{k+1} = ([\mathbf{H}_i^{k+1}]_s + l_i^{k+1} \mathbf{I}) w_i^{k+1} - \nabla f_i(w_i^{k+1})$

  Send  $\mathbf{S}_i^k$ ,  $l_i^{k+1} - l_i^k$ , and  $\xi_i^k$  to server

**for non-participating devices**  $i \notin S^k$  **do**

$z_i^{k+1} = z_i^k$ ,  $w_i^{k+1} = w_i^k$ ,  $\mathbf{L}_i^{k+1} = \mathbf{L}_i^k$ ,  $\mathbf{H}_i^{k+1} = \mathbf{H}_i^k$ ,  $l_i^{k+1} = l_i^k$ ,  $g_i^{k+1} = g_i^k$

**end**

**on server**

**if**  $\xi_i^k = 1$

$w_i^{k+1} = z_i^{k+1}$ , receive  $g_i^{k+1} - g_i^k$

**if**  $\xi_i^k = 0$

$w_i^{k+1} = w_i^k$ ,  $g_i^{k+1} - g_i^k = \alpha [\sum_{jl} (\mathbf{S}_i^k)_{jl} \mathbf{B}_i^{jl}]_s w_i^{k+1} + (l_i^{k+1} - l_i^k) w_i^{k+1}$

$g^{k+1} = g^k + \frac{1}{n} \sum_{i \in S^k} (g_i^{k+1} - g_i^k)$

$\mathbf{H}^{k+1} = \mathbf{H}^k + \frac{\alpha}{n} \sum_{i \in S^k} \sum_{jl} (\mathbf{S}_i^k)_{jl} \mathbf{B}_i^{jl}$

$l^{k+1} = l^k + \frac{1}{n} \sum_{i \in S^k} (l_i^{k+1} - l_i^k)$

## Convergence Result

Let  $\Phi_3^k \stackrel{\text{def}}{=} \mathcal{W}^k + \frac{2p}{A_M} \left(1 - \frac{\tau A_M}{n}\right) \mathcal{Z}^k$ , where  $\mathcal{Z}^k \stackrel{\text{def}}{=} \frac{1}{n} \sum_{i=1}^n \|z_i^k - x^*\|^2$ , for  $k \geq 0$ .

**Linear convergence of BL2:** Let Assumption 5 hold. Let Assumption 1 (i) or Assumption 2 (i) hold. Assume  $\|z_i^k - x^*\|^2 \leq \frac{A_M \mu^2}{(6H^2 + 24H_1^2) B_M}$  and  $\mathcal{H}^k \leq \frac{A_M \mu^2}{96N_B R^2 B_M}$  for all  $i \in [n]$  and  $k \geq 0$ . Then for  $k \geq 0$

$$\mathbb{E}[\Phi_3^k] \leq \left(1 - \frac{\tau \min\{p, A_M\}}{2n}\right)^k \Phi_3^0.$$

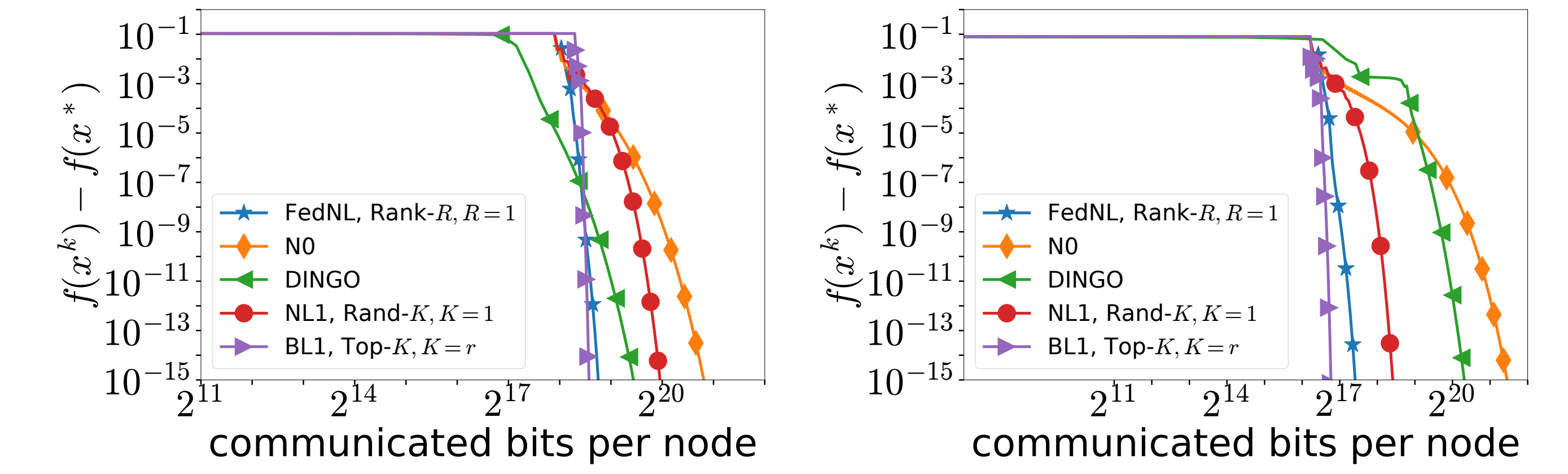
**Superlinear convergence of BL2:** Define  $\Phi_4^k \stackrel{\text{def}}{=} \mathcal{H}^k + \frac{4BM_M^2}{A_M} \|x^k - x^*\|^2$  for  $k \geq 0$ . Let  $\eta = 1$ ,  $\xi^k \equiv 1$ ,  $S^k \equiv [n]$ , and  $\mathcal{Q}_i^k(x) \equiv x$  for any  $x \in \mathbb{R}^d$  and  $k \geq 0$ . Let Assumption 5 hold. Let Assumption 3 (i) or Assumption 4 (i) hold. Assume  $\|z_i^k - x^*\|^2 \leq \frac{A_M \mu^2}{(6H^2 + 24H_1^2) B_M}$  and  $\mathcal{H}^k \leq \frac{A_M \mu^2}{96N_B R^2 B_M}$  for all  $i \in [n]$  and  $k \geq 0$ . Then we have

$$\mathbb{E}[\Phi_4^k] \leq \theta_2^k \Phi_4^0, \\ \mathbb{E} \left[ \frac{\|x^{k+1} - x^*\|^2}{\|x^k - x^*\|^2} \right] \leq \theta_2^k \left( \frac{A_M(3H^2 + 12H_1^2)}{16BM_M^2 \mu^2} + \frac{12N_B R^2}{\mu^2} \right) \Phi_4^0,$$

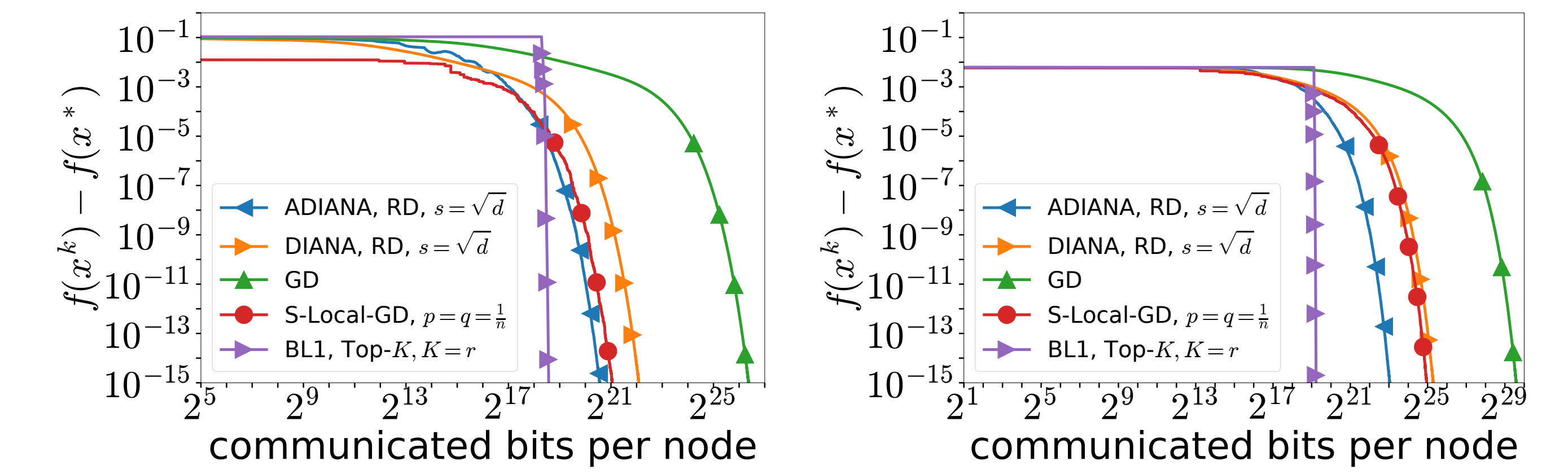
for  $k \geq 0$ , where  $\theta_2 \stackrel{\text{def}}{=} \left(1 - \frac{\min\{2A, A_M\}}{2}\right)$ .

## Numerical Results

1. BL1 vs N0 vs FedNL vs NL1 vs DINGO



2. BL1 vs DIANA vs ADIANA vs GD vs S-Local-GD



3. ECLK vs ADIANA

