

MARINA: Faster Non-Convex Distributed Learning with Compression

Eduard Gorbunov

MIPT, IITP

Konstantin Burlachenko
KAUST

Zhize Li
KAUST

Peter Richtárik
KAUST



Federated
Learning
One
World
Seminar



IITP RAS



March 10, 2021





Konstantin Burlachenko
PhD student
KAUST



Zhize Li
Research Scientist
KAUST

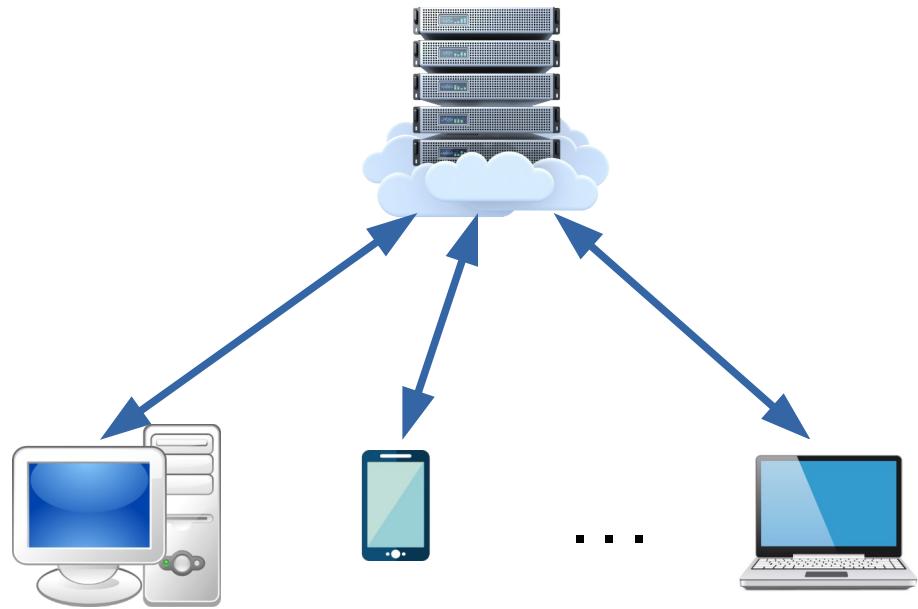


Peter Richtárik
Professor of Computer Science
KAUST

Outline

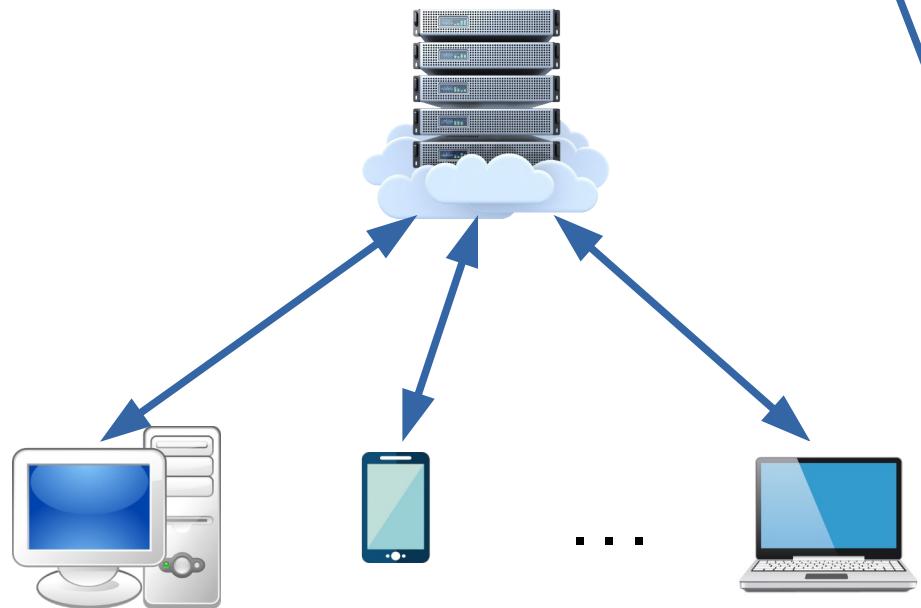
- 1 The problem
 - 2 Compressed communications
 - 3 Quantized Gradient Descent and DIANA
 - 4 MARINA
 - 5 MARINA and variance reduction
 - 6 MARINA and partial participation
 - 7 Experiments
-
- The diagram features three vertical blue curly braces on the right side of the list. The first brace covers sections 1 through 3, labeled '10 minutes' in blue text. The second brace covers sections 4 through 6, labeled '25 minutes' in blue text. The third brace covers section 7, labeled '5 minutes' in blue text.

1. The Problem



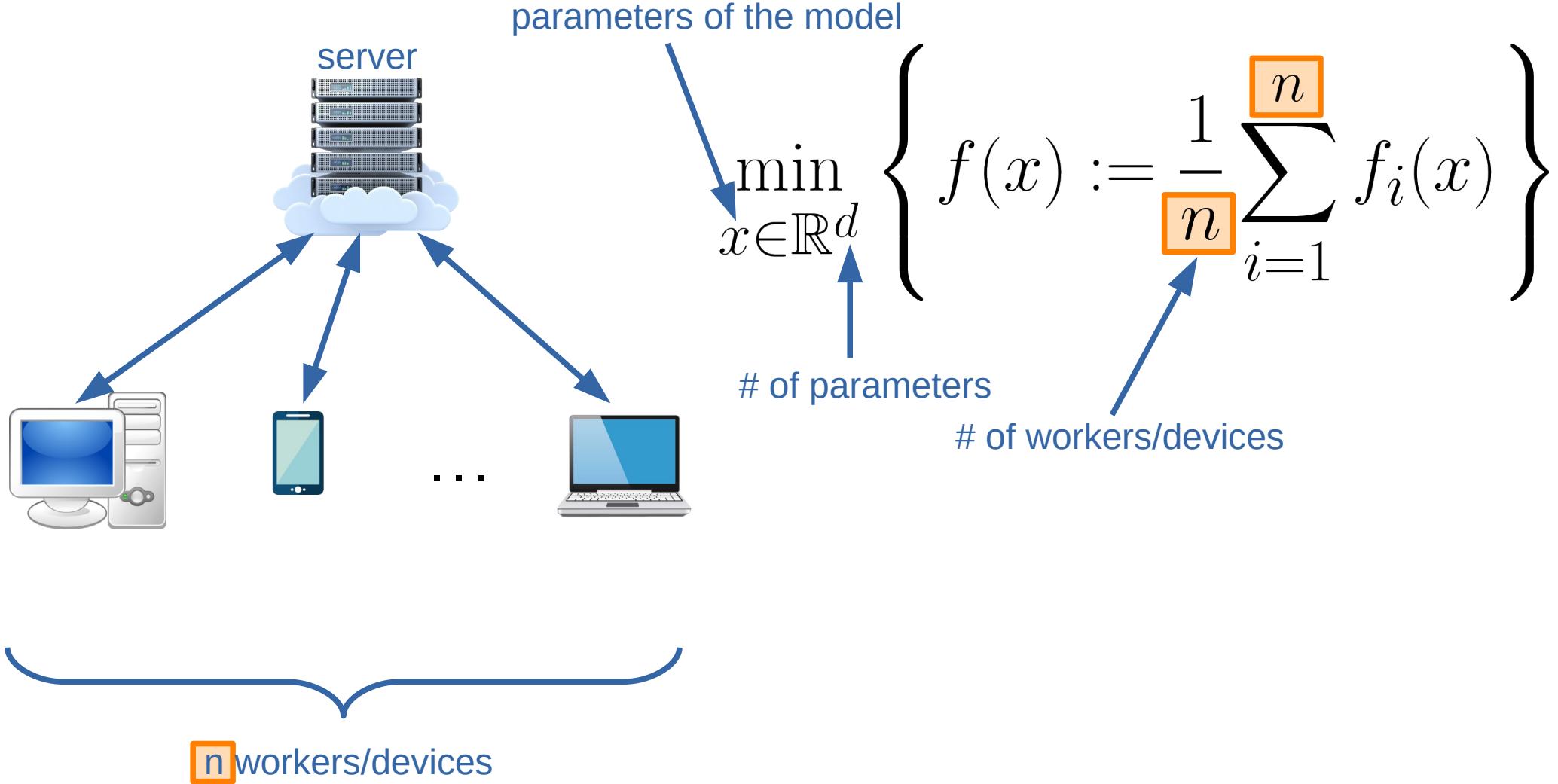
$$\min_{x \in \mathbb{R}^d} \left\{ f(x) := \frac{1}{n} \sum_{i=1}^n f_i(x) \right\}$$

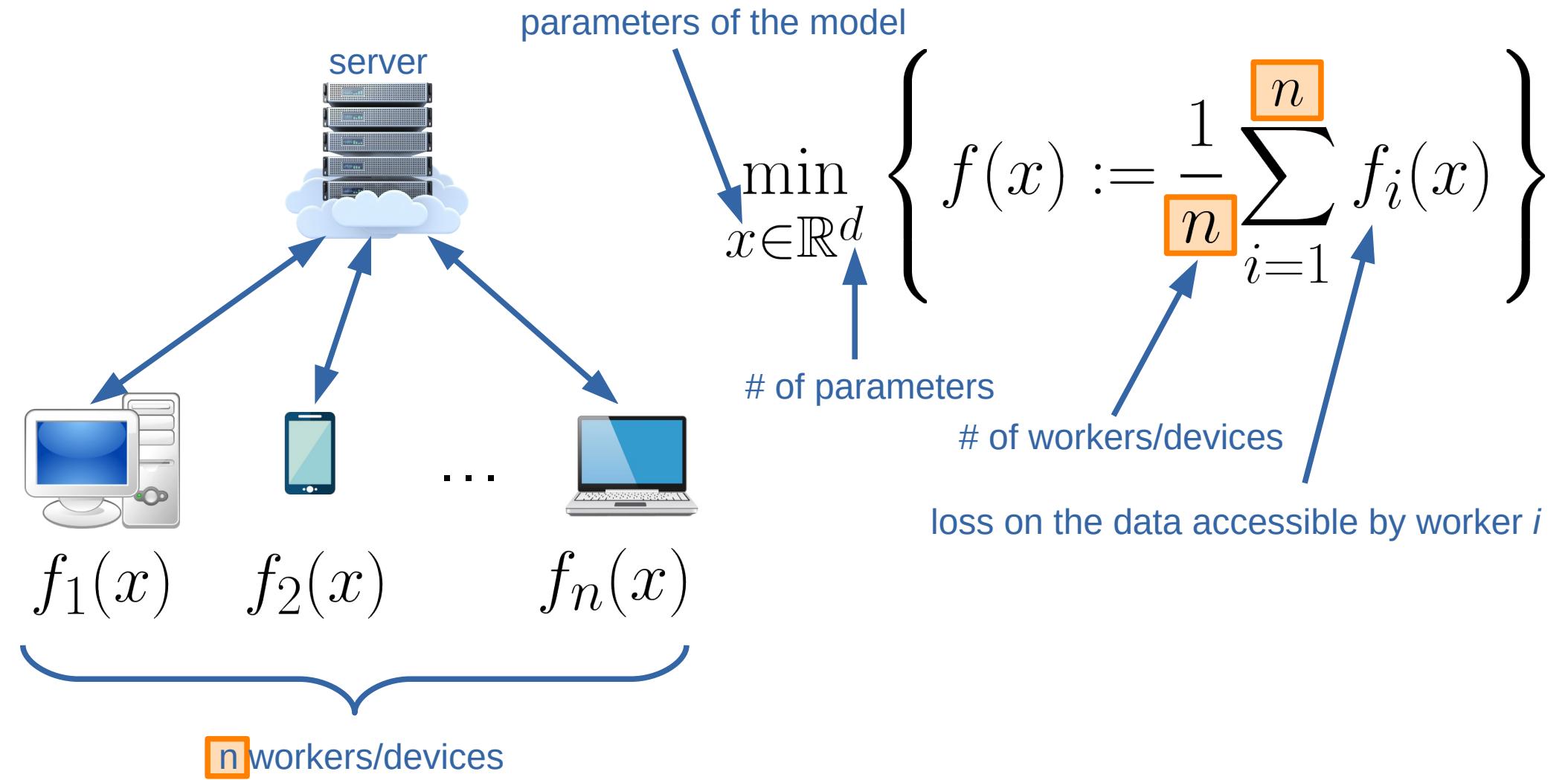
parameters of the model

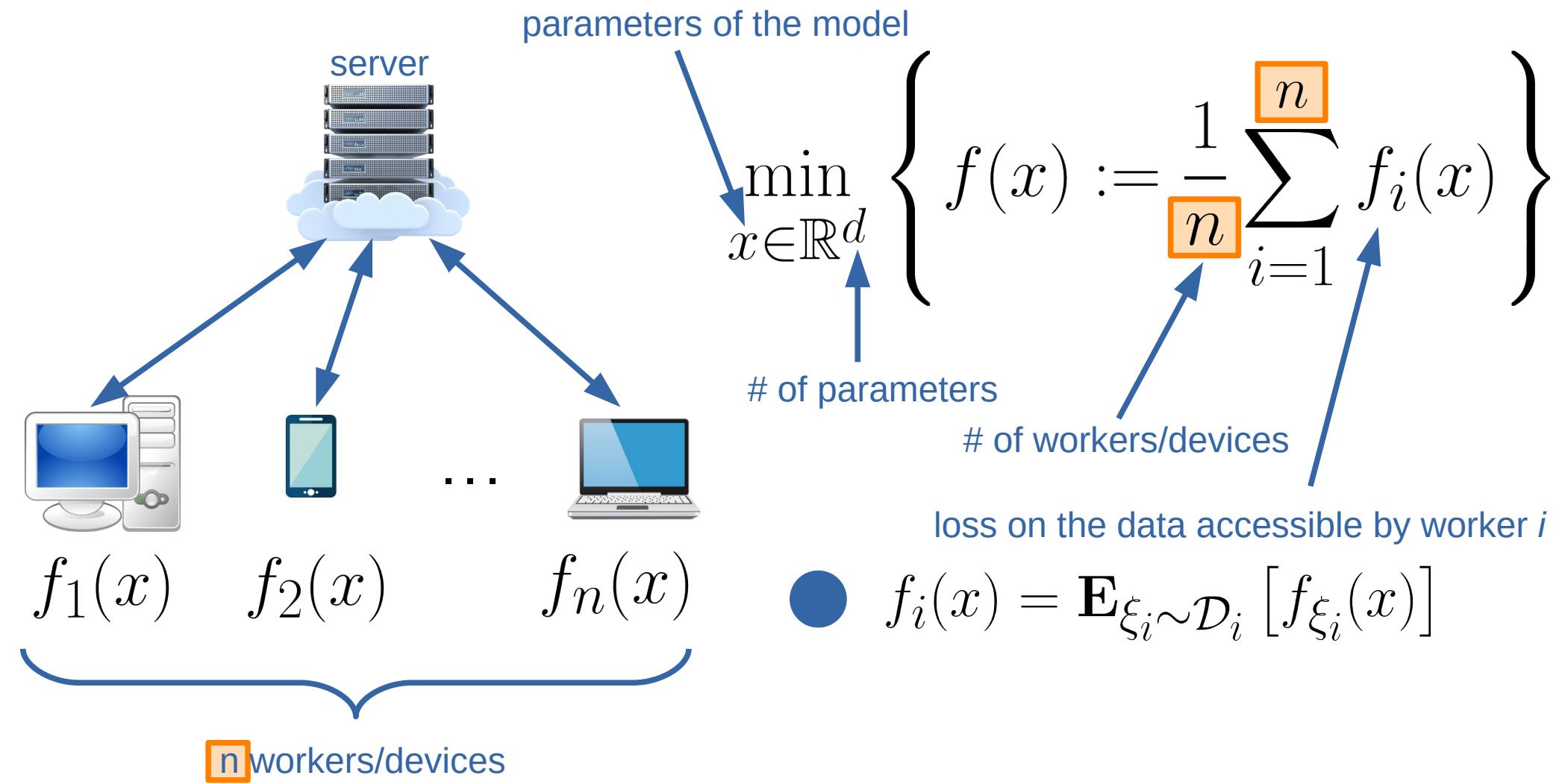


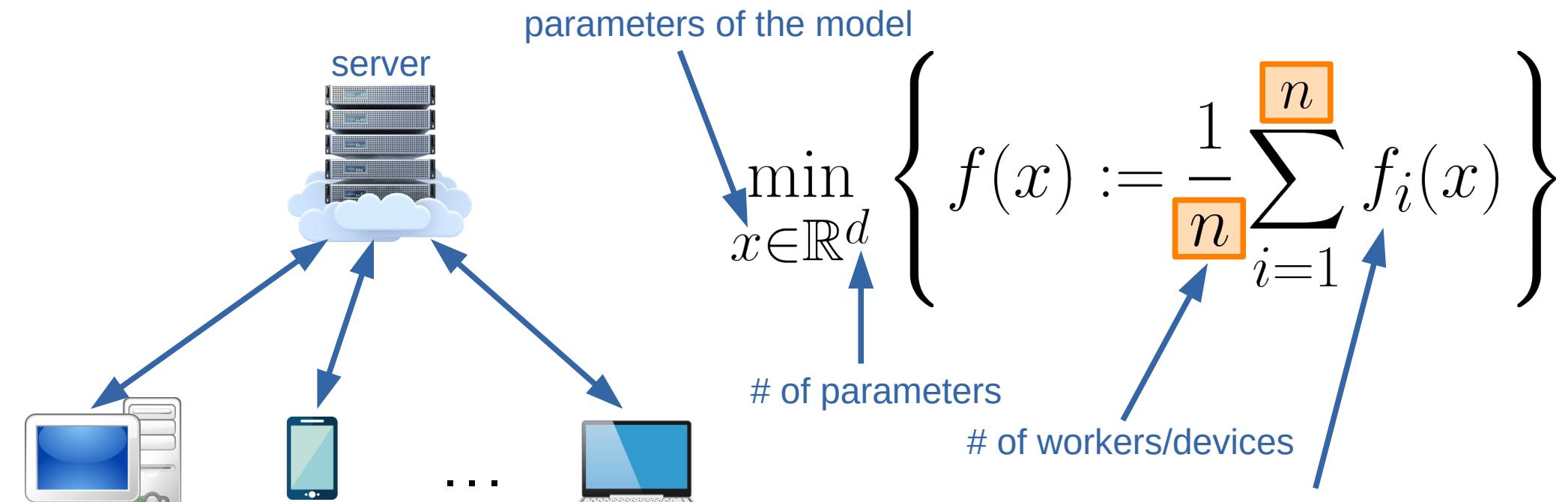
$$\min_{x \in \mathbb{R}^d} \left\{ f(x) := \frac{1}{n} \sum_{i=1}^n f_i(x) \right\}$$

of parameters









$$f_1(x) \quad f_2(x) \quad f_n(x)$$

n workers/devices

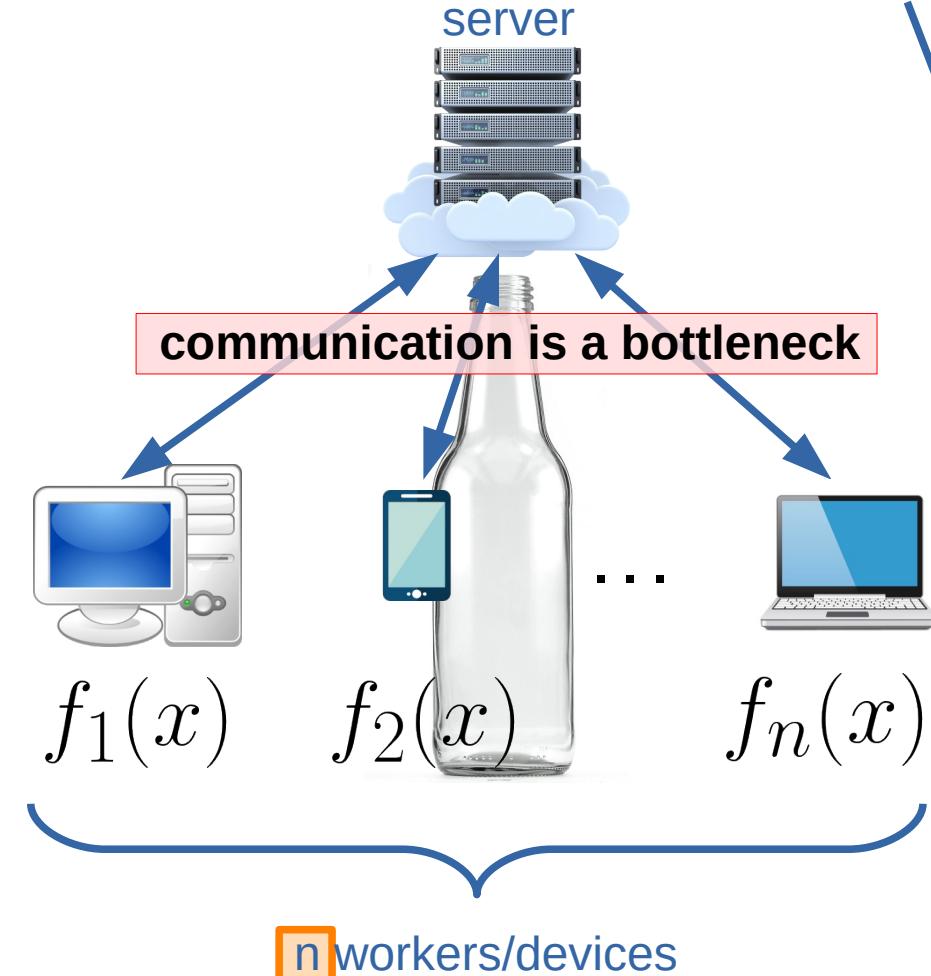


$$f_i(x) = \mathbf{E}_{\xi_i \sim \mathcal{D}_i} [f_{\xi_i}(x)]$$



$$f_i(x) = \frac{1}{m} \sum_{j=1}^m f_{ij}(x)$$

parameters of the model



$$\min_{x \in \mathbb{R}^d} \left\{ f(x) := \frac{1}{n} \sum_{i=1}^n f_i(x) \right\}$$

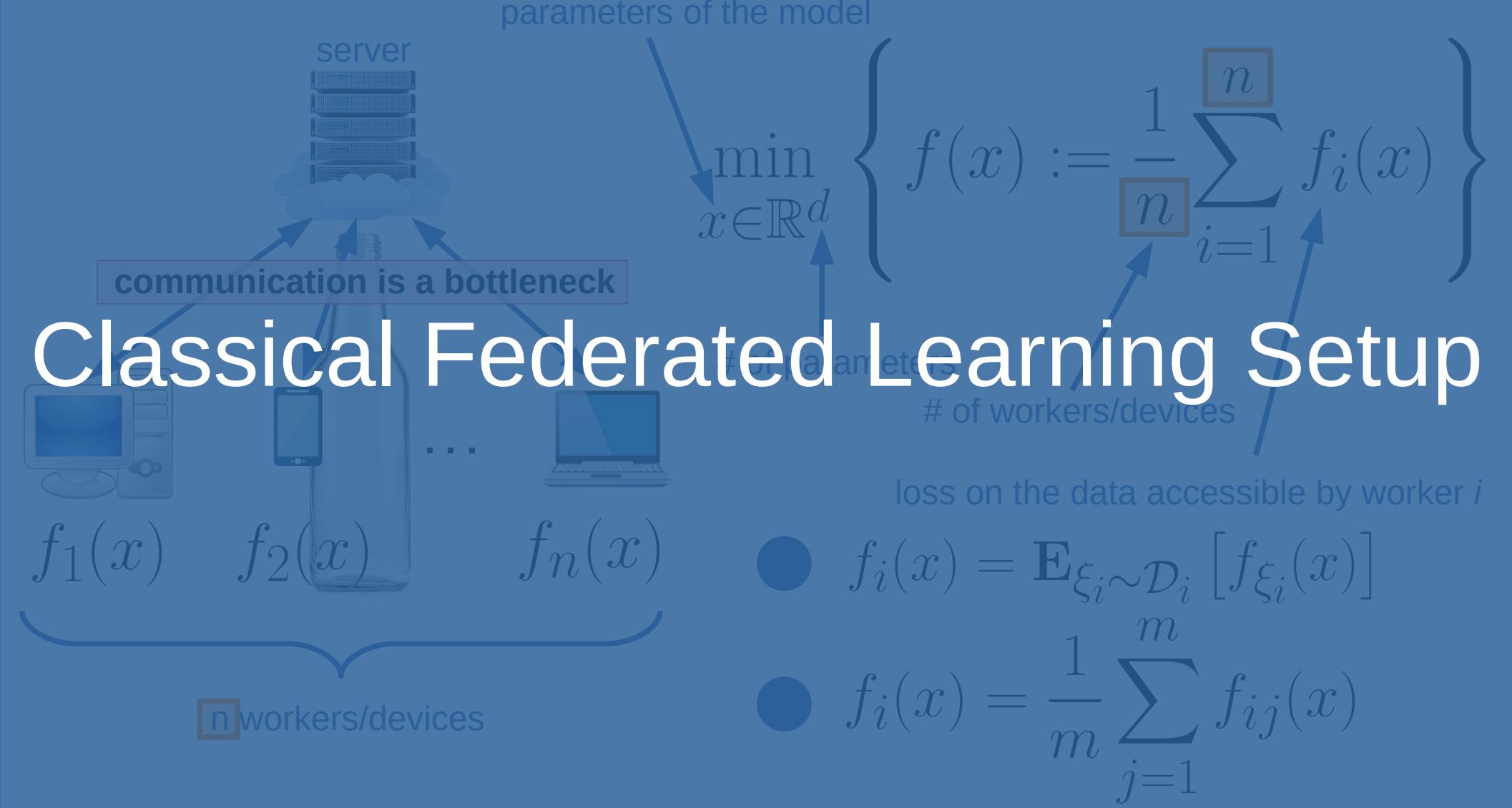
of parameters # of workers/devices

loss on the data accessible by worker i

● $f_i(x) = \mathbf{E}_{\xi_i \sim \mathcal{D}_i} [f_{\xi_i}(x)]$

● $f_i(x) = \frac{1}{m} \sum_{j=1}^m f_{ij}(x)$

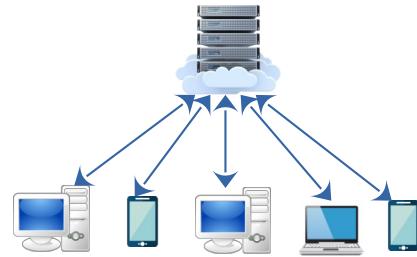
Classical Federated Learning Setup



How to Handle Communication Bottleneck?

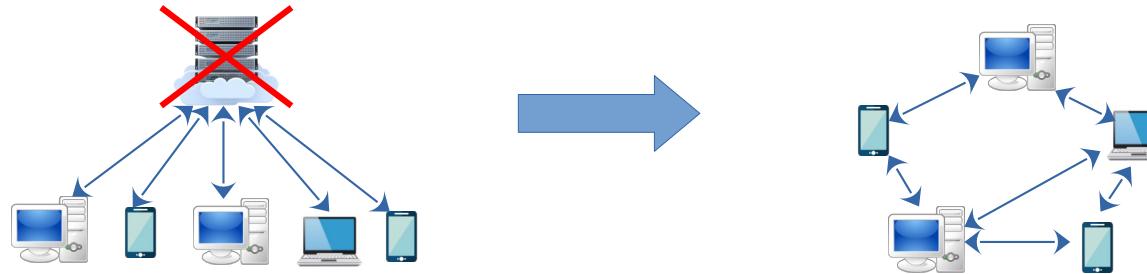
How to Handle Communication Bottleneck?

- Change the topology of the network



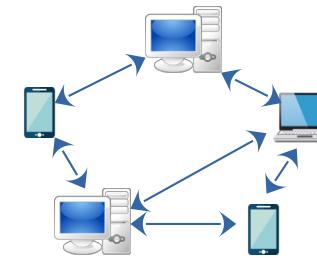
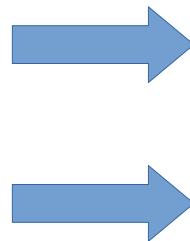
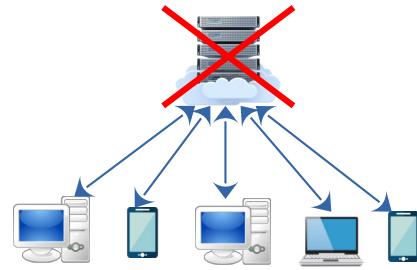
How to Handle Communication Bottleneck?

- Change the topology of the network → Decentralized optimization



How to Handle Communication Bottleneck?

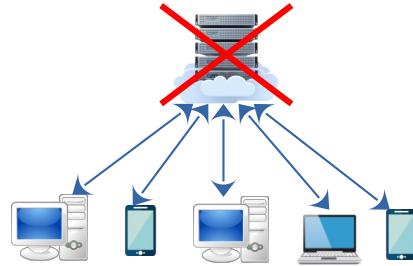
- Change the topology of the network → Decentralized optimization



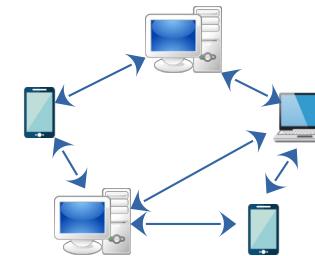
- Do more work on each worker in the hope of communicating less

How to Handle Communication Bottleneck?

- Change the topology of the network



Decentralized optimization

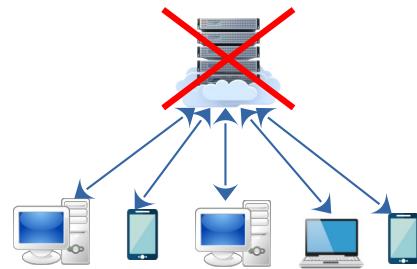


- Do more work on each worker in the hope of communicating less

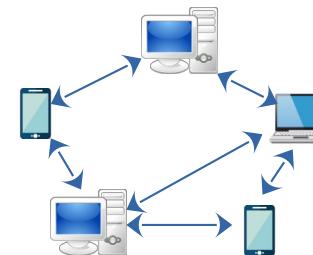
Local-SGD/Federated Averaging

How to Handle Communication Bottleneck?

- Change the topology of the network → Decentralized optimization



Decentralized optimization



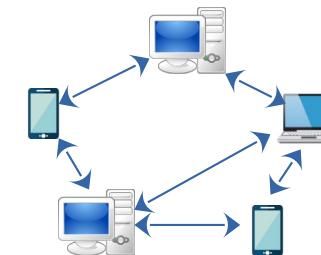
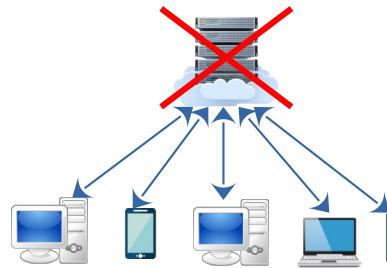
- Do more work on each worker in the hope of communicating less → Local-SGD/Federated Averaging
- Send less information to reduce the communication cost



Local-SGD/Federated Averaging

How to Handle Communication Bottleneck?

- Change the topology of the network → Decentralized optimization



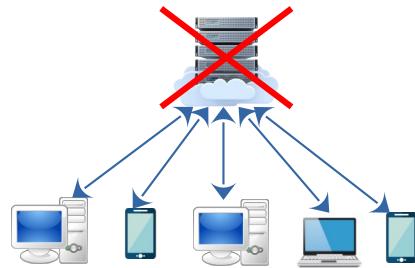
- Do more work on each worker in the hope of communicating less → Local-SGD/Federated Averaging
- Send less information to reduce the communication cost

Workers send dense vectors

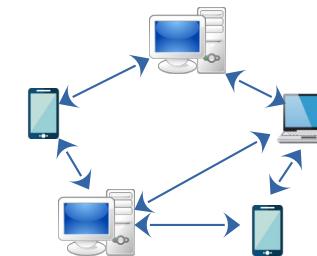
$$g = \begin{pmatrix} 1 \\ -15 \\ 0.2 \\ -7 \\ 10 \end{pmatrix}$$

How to Handle Communication Bottleneck?

- Change the topology of the network



Decentralized optimization



- Do more work on each worker in the hope of communicating less

Local-SGD/Federated Averaging

- Send less information to reduce the communication cost

Workers send dense vectors

$$g = \begin{pmatrix} 1 \\ -15 \\ 0.2 \\ -7 \\ 10 \end{pmatrix}$$



Workers send compressed/sparse vectors

$$\mathcal{Q}(g) = \frac{5}{2} \cdot \begin{pmatrix} 1 \\ 0 \\ 0 \\ -7 \\ 0 \end{pmatrix}$$

How to Handle Communication Bottleneck?

- Change the topology of the network → Decentralized optimization

We study this approach



- Do more work on each worker in the hope of communicating less → Local-SGD/Federated Averaging

- Send less information to reduce the communication cost

Workers send dense vectors

$$g = \begin{pmatrix} 1 \\ -15 \\ 0.2 \\ -7 \\ 10 \end{pmatrix}$$

Workers send compressed/sparse vectors

$$\mathcal{Q}(g) = \frac{5}{2} \cdot \begin{pmatrix} 1 \\ 0 \\ 0 \\ -7 \\ 0 \end{pmatrix}$$

Unbiased compression (quantization)

$$x \rightarrow Q(x)$$

Unbiased compression (quantization)

$$x \rightarrow Q(x) \quad \mathbb{E}[Q(x)] = x$$

Unbiased compression (quantization)

$$x \rightarrow Q(x) \quad \mathbb{E}[Q(x)] = x$$

$$\mathbb{E}\|Q(x) - x\|^2 \leq \omega \|x\|^2$$

Unbiased compression (quantization)

$$x \rightarrow Q(x) \quad \mathbb{E}[Q(x)] = x$$

$$\mathbb{E}\|Q(x) - x\|^2 \leq \omega \|x\|^2$$

Example: RandK (for K = 2)

$$\begin{pmatrix} 1 \\ -15 \\ 0.2 \\ -7 \\ 10 \end{pmatrix} \longrightarrow \text{blue arrow}$$

Unbiased compression (quantization)

$$x \rightarrow \mathcal{Q}(x) \quad \mathbb{E}[\mathcal{Q}(x)] = x$$

$$\mathbb{E}\|\mathcal{Q}(x) - x\|^2 \leq \omega\|x\|^2$$

Example: RandK (for K = 2)

$$\begin{pmatrix} 1 \\ -15 \\ 0.2 \\ -7 \\ 10 \end{pmatrix} \longrightarrow$$

Pick K = 2 components uniformly at random

Unbiased compression (quantization)

$$x \rightarrow \mathcal{Q}(x) \quad \mathbb{E}[\mathcal{Q}(x)] = x$$

$$\mathbb{E}\|\mathcal{Q}(x) - x\|^2 \leq \omega\|x\|^2$$

Example: RandK (for K = 2)

$$\begin{pmatrix} 1 \\ -15 \\ 0.2 \\ -7 \\ 10 \end{pmatrix} \xrightarrow{\text{blue arrow}} \frac{5}{2} \cdot \begin{pmatrix} 1 \\ 0 \\ 0 \\ -7 \\ 0 \end{pmatrix}$$

Pick K = 2 components uniformly at random

Unbiased compression (quantization)

$$x \rightarrow \mathcal{Q}(x) \quad \mathbb{E}[\mathcal{Q}(x)] = x$$

$$\mathbb{E}\|\mathcal{Q}(x) - x\|^2 \leq \omega\|x\|^2$$

Example: RandK (for K = 2)

$$\begin{pmatrix} 1 \\ -15 \\ 0.2 \\ -7 \\ 10 \end{pmatrix} \xrightarrow{\text{for unbiasedness}} \frac{5}{2} \cdot \begin{pmatrix} 1 \\ 0 \\ 0 \\ -7 \\ 0 \end{pmatrix}$$

Pick K = 2 components uniformly at random

Unbiased compression (quantization)

$$x \rightarrow Q(x) \quad \mathbb{E}[Q(x)] = x$$

$$\mathbb{E}\|Q(x) - x\|^2 \leq \omega \|x\|^2$$

Example: RandK (for K = 2)

$$\begin{pmatrix} 1 \\ -15 \\ 0.2 \\ -7 \\ 10 \end{pmatrix} \xrightarrow{\text{for unbiasedness}} \begin{matrix} 5 \\ 2 \end{matrix} \cdot \begin{pmatrix} 1 \\ 0 \\ 0 \\ -7 \\ 0 \end{pmatrix}$$

Pick K = 2 components uniformly at random

Unbiased compression (quantization)

$$x \rightarrow \mathcal{Q}(x) \quad \mathbb{E}[\mathcal{Q}(x)] = x$$

$$\mathbb{E}\|\mathcal{Q}(x) - x\|^2 \leq \omega\|x\|^2$$

Example: RandK (for K = 2)

$$\begin{pmatrix} 1 \\ -15 \\ 0.2 \\ -7 \\ 10 \end{pmatrix} \xrightarrow{\text{for unbiasedness}} 5 \cdot \begin{pmatrix} 1 \\ 0 \\ 0 \\ -7 \\ 0 \end{pmatrix}$$

$\omega = \frac{d}{K} - 1$

Pick K = 2 components uniformly at random

Unbiased compression (quantization)

$$x \rightarrow Q(x) \quad \mathbb{E}[Q(x)] = x$$

$$\mathbb{E}\|Q(x) - x\|^2 \leq \omega \|x\|^2$$

Example: RandK (for K = 2)

$$d = 5 \left\{ \begin{pmatrix} 1 \\ -15 \\ 0.2 \\ -7 \\ 10 \end{pmatrix} \right. \xrightarrow{\text{for unbiasedness}} \left. \begin{matrix} 5 \\ 2 \end{matrix} \right\} \cdot \begin{pmatrix} 1 \\ 0 \\ 0 \\ -7 \\ 0 \end{pmatrix}$$

$$\omega = \frac{d}{K} - 1$$

Pick K = 2 components uniformly at random

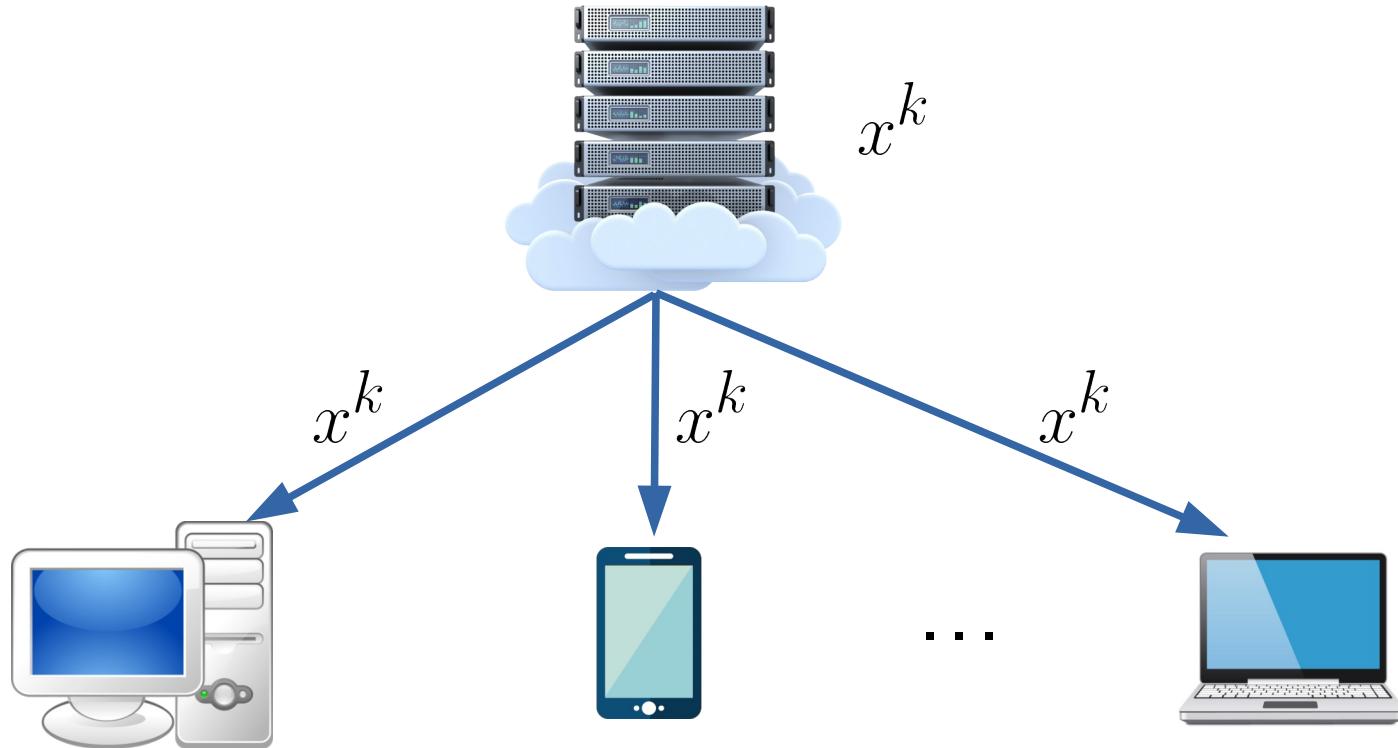
2. Quantized Gradient Descent (QGD)



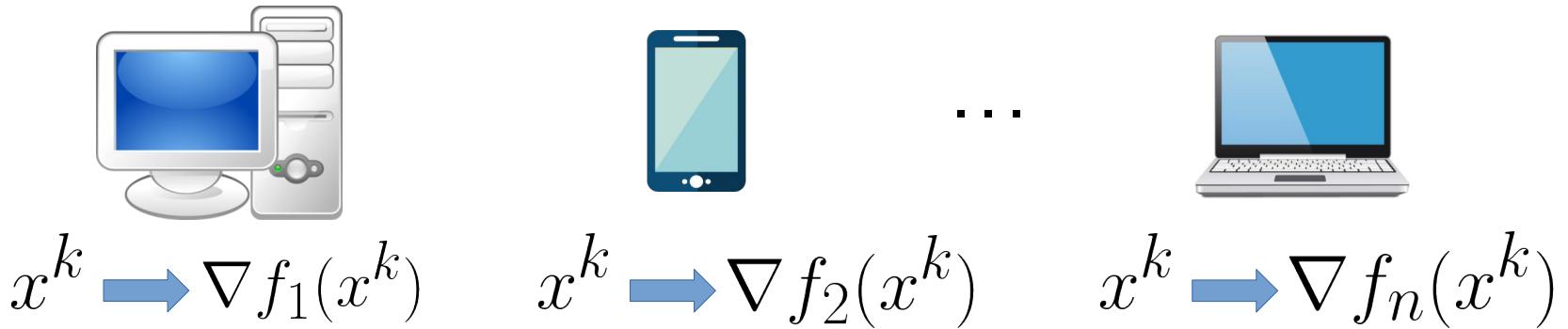
Alistarh, Dan, Demjan Grubic, Jerry Li, Ryota Tomioka, and Milan Vojnovic. "**QSGD: Communication-efficient SGD via gradient quantization and encoding.**" *In Advances in Neural Information Processing Systems*, pp. 1709-1720. 2017.

1

Server broadcasts the parameters



- 1 Server broadcasts the parameters
- 2 Devices compute the gradients



1 Server broadcasts the parameters

2 Devices compute the gradients

3 Devices quantize the gradients



x^k



$$x^k \rightarrow \nabla f_1(x^k)$$



...



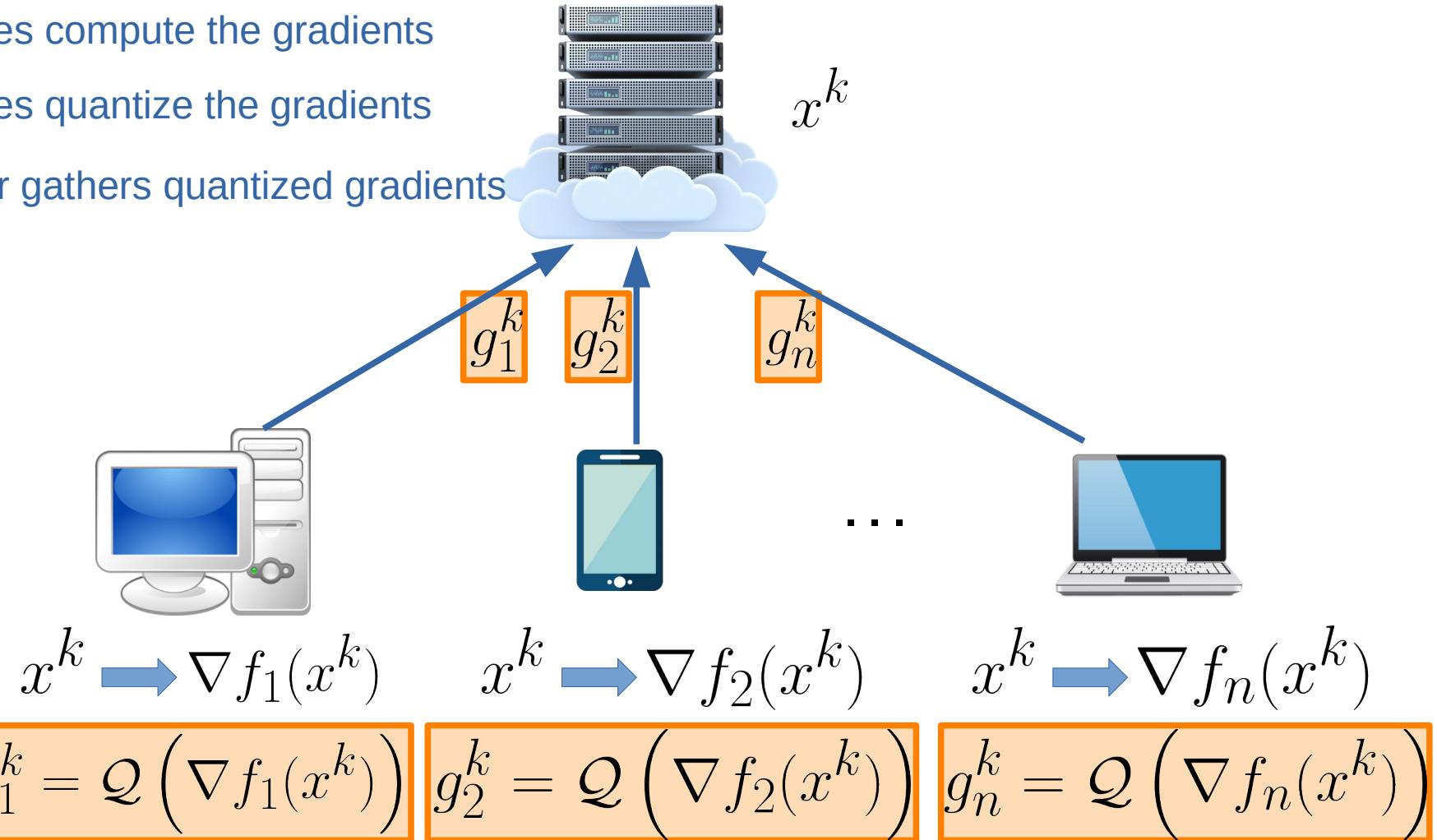
$$x^k \rightarrow \nabla f_n(x^k)$$

$$g_1^k = \mathcal{Q}(\nabla f_1(x^k))$$

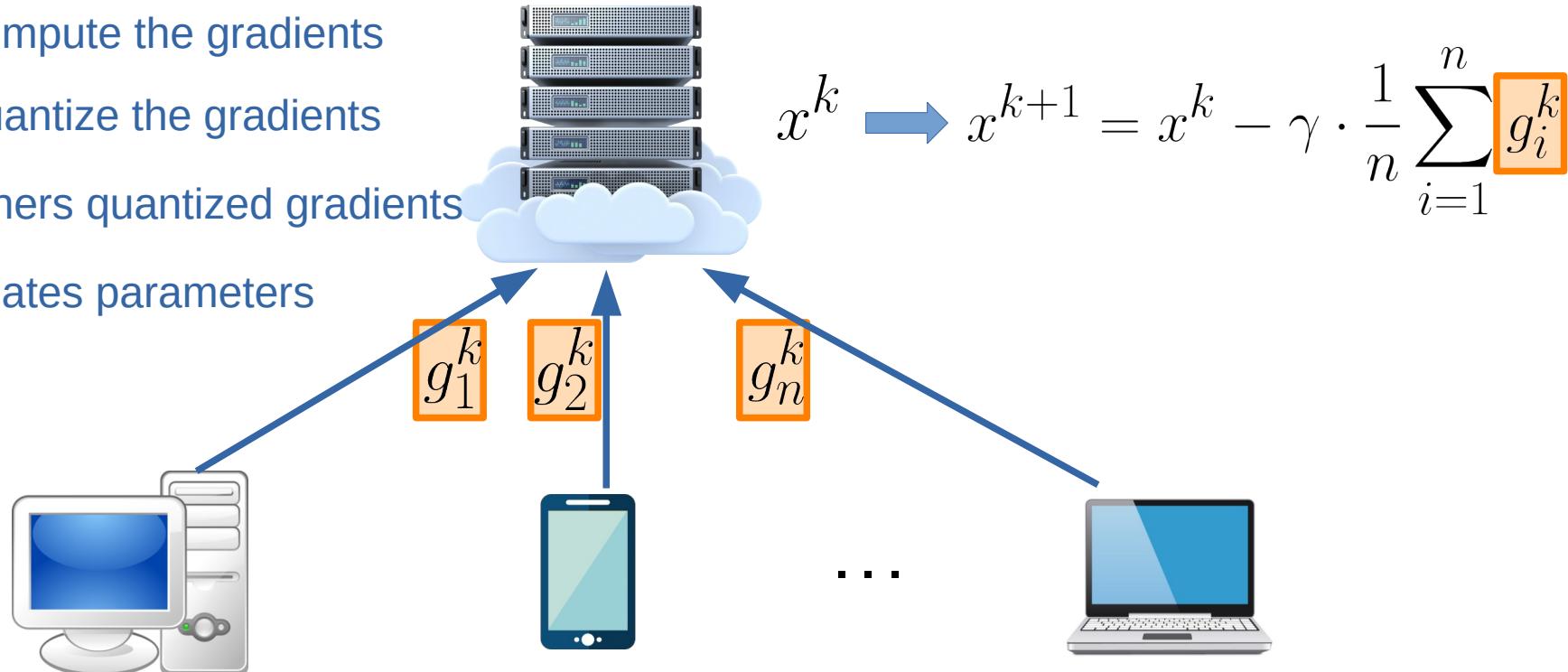
$$g_2^k = \mathcal{Q}(\nabla f_2(x^k))$$

$$g_n^k = \mathcal{Q}(\nabla f_n(x^k))$$

- 1 Server broadcasts the parameters
- 2 Devices compute the gradients
- 3 Devices quantize the gradients
- 4 Server gathers quantized gradients



- 1 Server broadcasts the parameters
- 2 Devices compute the gradients
- 3 Devices quantize the gradients
- 4 Server gathers quantized gradients
- 5 Server updates parameters



$$x^k \rightarrow \nabla f_1(x^k)$$

$$x^k \rightarrow \nabla f_2(x^k)$$

$$x^k \rightarrow \nabla f_n(x^k)$$

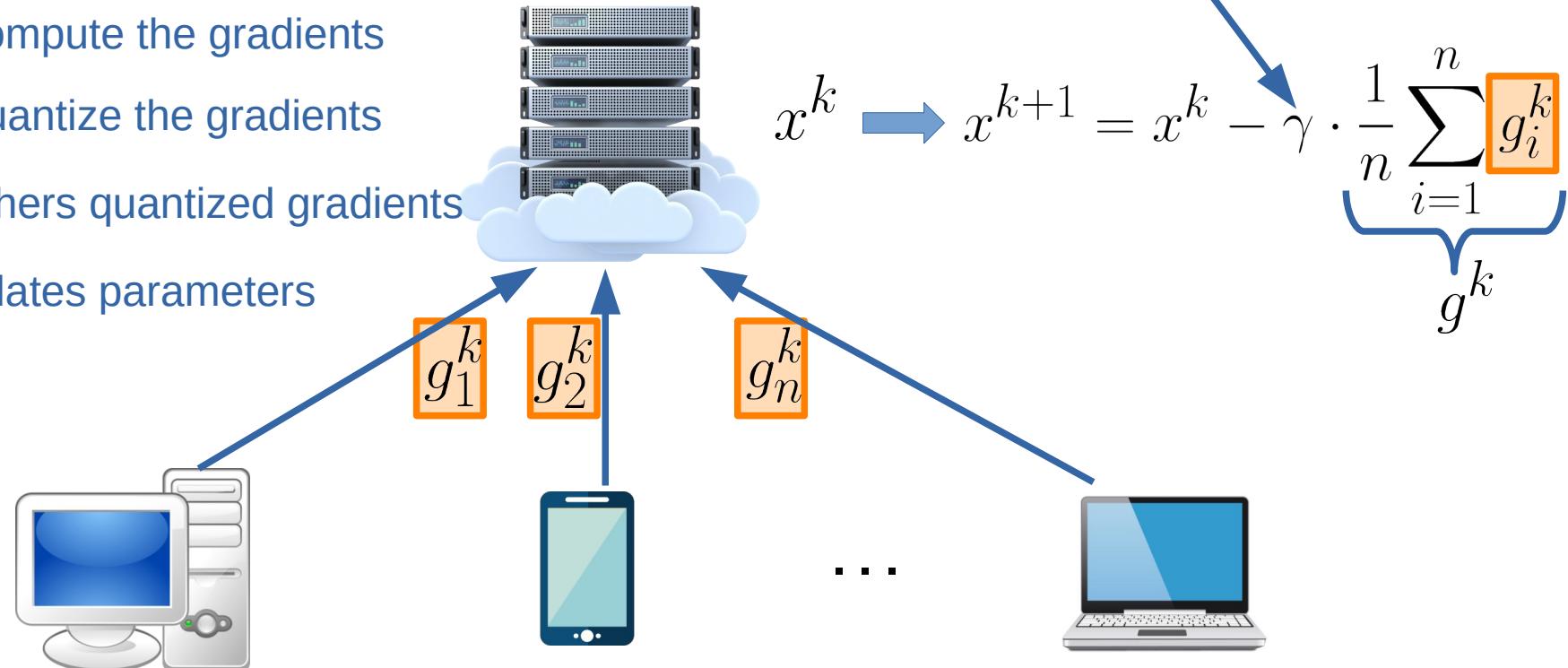
$$g_1^k = \mathcal{Q} \left(\nabla f_1(x^k) \right)$$

$$g_2^k = \mathcal{Q} \left(\nabla f_2(x^k) \right)$$

$$g_n^k = \mathcal{Q} \left(\nabla f_n(x^k) \right)$$

$$x^k \rightarrow x^{k+1} = x^k - \gamma \cdot \frac{1}{n} \sum_{i=1}^n g_i^k$$

- 1 Server broadcasts the parameters
- 2 Devices compute the gradients
- 3 Devices quantize the gradients
- 4 Server gathers quantized gradients
- 5 Server updates parameters



$$x^k \rightarrow \nabla f_1(x^k)$$

$$x^k \rightarrow \nabla f_2(x^k)$$

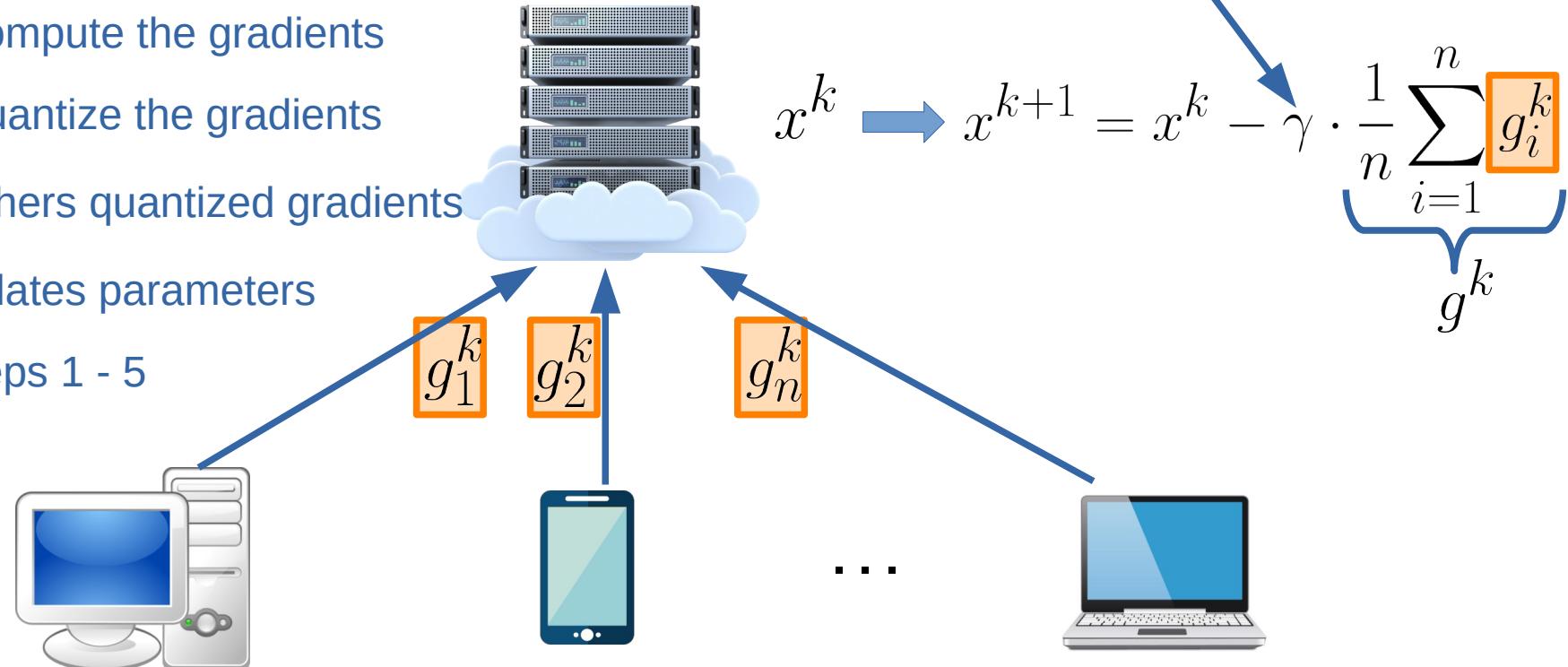
$$x^k \rightarrow \nabla f_n(x^k)$$

$$g_1^k = Q(\nabla f_1(x^k))$$

$$g_2^k = Q(\nabla f_2(x^k))$$

$$g_n^k = Q(\nabla f_n(x^k))$$

- 1 Server broadcasts the parameters
- 2 Devices compute the gradients
- 3 Devices quantize the gradients
- 4 Server gathers quantized gradients
- 5 Server updates parameters
- 6 Repeat steps 1 - 5



$$x^k \rightarrow \nabla f_1(x^k)$$

$$x^k \rightarrow \nabla f_2(x^k)$$

$$x^k \rightarrow \nabla f_n(x^k)$$

$$g_1^k = Q(\nabla f_1(x^k))$$

$$g_2^k = Q(\nabla f_2(x^k))$$

$$g_n^k = Q(\nabla f_n(x^k))$$

Assumptions

- 1 Uniform lower bound:

Assumptions

- 1 Uniform lower bound:

$$\exists f_* \in \mathbb{R} : \forall x \in \mathbb{R}^d \quad f(x) \geq f_*$$

Assumptions

- 1 Uniform lower bound:

$$\exists f_* \in \mathbb{R} : \forall x \in \mathbb{R}^d \quad f(x) \geq f_*$$

- 2 Smoothness:

Assumptions

1 Uniform lower bound:

$$\exists f_* \in \mathbb{R} : \forall x \in \mathbb{R}^d \quad f(x) \geq f_*$$

2 Smoothness:

$$\|\nabla f_i(x) - \nabla f_i(y)\| \leq L_i \|x - y\|$$

Complexity Bound for QGD



Khaled, Ahmed, and Peter Richtárik. "Better theory for SGD in the nonconvex world." arXiv preprint arXiv:2002.03329 (2020).

QGD finds such \hat{x} that $\mathbb{E} \left[\|\nabla f(\hat{x})\|^2 \right] \leq \varepsilon^2$ after

Complexity Bound for QGD



Khaled, Ahmed, and Peter Richtárik. "Better theory for SGD in the nonconvex world." arXiv preprint arXiv:2002.03329 (2020).

QGD finds such \hat{x} that $\mathbb{E} \left[\|\nabla f(\hat{x})\|^2 \right] \leq \varepsilon^2$ after

$$\mathcal{O} \left(\frac{\Delta_0}{\varepsilon^2} + \frac{(1 + \omega)\Delta_0^2}{\varepsilon^4 n} + \frac{(1 + \omega)\Delta_0\Delta_f^*}{\varepsilon^4 n} \right)$$

communication
rounds

Complexity Bound for QGD



Khaled, Ahmed, and Peter Richtárik. "Better theory for SGD in the nonconvex world." arXiv preprint arXiv:2002.03329 (2020).

QGD finds such \hat{x} that $\mathbb{E} \left[\|\nabla f(\hat{x})\|^2 \right] \leq \varepsilon^2$ after

Hides
numerical
factors and
smoothness
constants

$$\mathcal{O} \left(\frac{\Delta_0}{\varepsilon^2} + \frac{(1 + \omega)\Delta_0^2}{\varepsilon^4 n} + \frac{(1 + \omega)\Delta_0\Delta_f^*}{\varepsilon^4 n} \right)$$

communication
rounds

Complexity Bound for QGD



Khaled, Ahmed, and Peter Richtárik. "Better theory for SGD in the nonconvex world." arXiv preprint arXiv:2002.03329 (2020).

QGD finds such \hat{x} that $\mathbb{E} \left[\|\nabla f(\hat{x})\|^2 \right] \leq \boxed{\varepsilon^2}$ after

Hides numerical factors and smoothness constants

$$\mathcal{O} \left(\frac{\Delta_0}{\boxed{\varepsilon^2}} + \frac{(1 + \omega)\Delta_0^2}{\boxed{\varepsilon^4} n} + \frac{(1 + \omega)\Delta_0\Delta_f^*}{\boxed{\varepsilon^4} n} \right)$$

communication rounds

Complexity Bound for QGD



Khaled, Ahmed, and Peter Richtárik. "Better theory for SGD in the nonconvex world." arXiv preprint arXiv:2002.03329 (2020).

QGD finds such \hat{x} that $\mathbb{E} \left[\|\nabla f(\hat{x})\|^2 \right] \leq \boxed{\varepsilon^2}$ after

Hides numerical factors and smoothness constants

$$\mathcal{O} \left(\frac{\Delta_0}{\boxed{\varepsilon^2}} + \frac{(1 + \boxed{\omega})\Delta_0^2}{\boxed{\varepsilon^4} n} + \frac{(1 + \boxed{\omega})\Delta_0\Delta_f^*}{\boxed{\varepsilon^4} n} \right)$$

communication rounds

$$\mathbb{E} \| \mathcal{Q}(x) - x \|^2 \leq \omega \| x \|^2$$

Complexity Bound for QGD



Khaled, Ahmed, and Peter Richtárik. "Better theory for SGD in the nonconvex world." arXiv preprint arXiv:2002.03329 (2020).

QGD finds such \hat{x} that $\mathbb{E} \left[\|\nabla f(\hat{x})\|^2 \right] \leq \varepsilon^2$ after

Hides numerical factors and smoothness constants $\rightarrow \mathcal{O} \left(\frac{\Delta_0}{\varepsilon^2} + \frac{(1 + \omega) \Delta_0^2}{\varepsilon^4 n} + \frac{(1 + \omega) \Delta_0 \Delta_f^*}{\varepsilon^4 n} \right)$ communication rounds

$$\mathcal{O} \left(\frac{\Delta_0}{\varepsilon^2} + \frac{(1 + \omega) \Delta_0^2}{\varepsilon^4 n} + \frac{(1 + \omega) \Delta_0 \Delta_f^*}{\varepsilon^4 n} \right)$$

$$\mathbb{E} \| \mathcal{Q}(x) - x \|^2 \leq \omega \| x \|^2$$

$$\Delta_0 = f(x^0) - f_*$$

Complexity Bound for QGD



Khaled, Ahmed, and Peter Richtárik. "Better theory for SGD in the nonconvex world." arXiv preprint arXiv:2002.03329 (2020).

QGD finds such \hat{x} that $\mathbb{E} \left[\|\nabla f(\hat{x})\|^2 \right] \leq \varepsilon^2$ after

Hides numerical factors and smoothness constants

$$\mathcal{O} \left(\frac{\Delta_0}{\varepsilon^2} + \frac{(1 + \omega) \Delta_0^2}{\varepsilon^4 n} + \frac{(1 + \omega) \Delta_0 \Delta_f^*}{\varepsilon^4 n} \right)$$

communication rounds

$$\mathbb{E} \| \mathcal{Q}(x) - x \|^2 \leq \omega \| x \|^2$$

$$\Delta_0 = f(x^0) - f_*$$

$$\Delta_f^* = f_* - \frac{1}{n} \sum_{i=1}^n f_{i,*}$$

3. DIANA



Mishchenko, Konstantin, Eduard Gorbunov, Martin Takáč, and Peter Richtárik.
"Distributed learning with compressed gradient differences." arXiv preprint
arXiv:1901.09269 (2019).



Horváth, Samuel, Dmitry Kovalev, Konstantin Mishchenko, Sebastian Stich, and Peter
Richtárik. **"Stochastic distributed learning with gradient quantization and variance
reduction."** arXiv preprint arXiv:1904.05115 (2019).

$$x^{k+1} = x^k - \gamma \cdot \frac{1}{n} \sum_{i=1}^n g_i^k$$

$$x^{k+1} = x^k - \gamma \cdot \frac{1}{n} \sum_{i=1}^n g_i^k$$

QGD: $g_i^k = \mathcal{Q}\left(\nabla f_i(x^k)\right)$

$$x^{k+1} = x^k - \gamma \cdot \frac{1}{n} \sum_{i=1}^n g_i^k$$

QGD: $g_i^k = \mathcal{Q}\left(\nabla f_i(x^k)\right)$

DIANA: $g_i^k = h_i^k + \mathcal{Q}\left(\nabla f_i(x^k) - h_i^k\right)$

$$x^{k+1} = x^k - \gamma \cdot \frac{1}{n} \sum_{i=1}^n g_i^k$$

QGD: $g_i^k = \mathcal{Q}\left(\nabla f_i(x^k)\right)$

DIANA: $g_i^k = h_i^k + \mathcal{Q}\left(\nabla f_i(x^k) - h_i^k\right)$



learnable local shifts

$$h_i^{k+1} = h_i^k + \alpha \mathcal{Q}\left(\nabla f_i(x^k) - h_i^k\right)$$

$$x^{k+1} = x^k - \gamma \cdot \frac{1}{n} \sum_{i=1}^n g_i^k$$

QGD: $g_i^k = \mathcal{Q}(\nabla f_i(x^k))$ vectors that devices have to send

DIANA: $g_i^k = h_i^k + \mathcal{Q}(\nabla f_i(x^k) - h_i^k)$ learnable local shifts

$$h_i^{k+1} = h_i^k + \alpha \mathcal{Q}(\nabla f_i(x^k) - h_i^k)$$

Complexity Bound for DIANA

DIANA finds such \hat{x} that $\mathbb{E} \left[\|\nabla f(\hat{x})\|^2 \right] \leq \varepsilon^2$ after

Complexity Bound for DIANA

DIANA finds such \hat{x} that $\mathbb{E} \left[\|\nabla f(\hat{x})\|^2 \right] \leq \varepsilon^2$ after

$$\mathcal{O} \left(\frac{\Delta_0 \left(1 + (1 + \omega) \sqrt{\omega/n} \right)}{\varepsilon^2} \right)$$

communication
rounds

Complexity Bound for DIANA

DIANA finds such \hat{x} that $\mathbb{E} \left[\|\nabla f(\hat{x})\|^2 \right] \leq \varepsilon^2$ after

Hides
numerical
factors and
smoothness
constants

$$\xrightarrow{\text{Hides numerical factors and smoothness constants}} O \left(\frac{\Delta_0 \left(1 + (1 + \omega) \sqrt{\omega/n} \right)}{\varepsilon^2} \right)$$

communication
rounds

Complexity Bound for DIANA

DIANA finds such \hat{x} that $\mathbb{E} \left[\|\nabla f(\hat{x})\|^2 \right] \leq \boxed{\varepsilon^2}$ after

Hides
numerical
factors and
smoothness
constants

$$\xrightarrow{\text{Hides numerical factors and smoothness constants}} O \left(\frac{\Delta_0 \left(1 + (1 + \omega) \sqrt{\omega/n} \right)}{\boxed{\varepsilon^2}} \right)$$

communication
rounds

Complexity Bound for DIANA

DIANA finds such \hat{x} that $\mathbb{E} \left[\|\nabla f(\hat{x})\|^2 \right] \leq \boxed{\varepsilon^2}$ after

Hides
numerical
factors and
smoothness
constants

$$\xrightarrow{\text{Hides numerical factors and smoothness constants}} \mathcal{O} \left(\frac{\Delta_0 \left(1 + (1 + \boxed{\omega}) \sqrt{\boxed{\omega/n}} \right)}{\boxed{\varepsilon^2}} \right)$$

communication
rounds

$$\boxed{\mathbb{E} \| \mathcal{Q}(x) - x \|^2 \leq \omega \| x \|^2}$$

Complexity Bound for DIANA

DIANA finds such \hat{x} that $\mathbb{E} \left[\|\nabla f(\hat{x})\|^2 \right] \leq \varepsilon^2$ after

Hides numerical factors and smoothness constants

$$\mathcal{O} \left(\frac{\Delta_0 \left(1 + (1 + \omega) \sqrt{\omega/n} \right)}{\varepsilon^2} \right)$$

communication rounds

$$\mathbb{E} \|\mathcal{Q}(x) - x\|^2 \leq \omega \|x\|^2$$

$$\Delta_0 = f(x^0) - f_*$$

Complexity Bounds for DIANA and QGD

QGD: $\mathcal{O} \left(\frac{\Delta_0}{\varepsilon^2} + \frac{(1+\omega)\Delta_0^2}{\varepsilon^4 n} + \frac{(1+\omega)\Delta_0\Delta_f^*}{\varepsilon^4 n} \right)$

DIANA: $\mathcal{O} \left(\frac{\Delta_0 \left(1 + (1+\omega)\sqrt{\omega/n} \right)}{\varepsilon^2} \right)$

Complexity Bounds for DIANA and QGD

QGD: $\mathcal{O} \left(\frac{\Delta_0}{\varepsilon^2} + \frac{(1+\omega)\Delta_0^2}{\varepsilon^4 n} + \frac{(1+\omega)\Delta_0\Delta_f^*}{\varepsilon^4 n} \right)$

DIANA: $\mathcal{O} \left(\frac{\Delta_0 \left(1 + (1+\omega)\sqrt{\omega/n} \right)}{\varepsilon^2} \right)$

Complexity Bounds for DIANA and QGD

QGD: $\mathcal{O} \left(\frac{\Delta_0}{\varepsilon^2} + \frac{(1 + \boxed{\omega})\Delta_0^2}{\varepsilon^4 n} + \frac{(1 + \boxed{\omega})\Delta_0 \Delta_f^*}{\varepsilon^4 n} \right)$

DIANA: $\mathcal{O} \left(\frac{\Delta_0 \left(1 + (1 + \boxed{\omega}) \sqrt{\boxed{\omega}/n} \right)}{\varepsilon^2} \right)$

Complexity Bound for DIANA

QGD: $\mathcal{O} \left(\frac{\Delta_0}{\varepsilon^2} + \frac{(1 + \omega)\Delta_0^2}{\varepsilon^4 n} + \frac{(1 + \omega)\Delta_0\Delta_f^*}{\varepsilon^4 n} \right)$

Is it possible to get better rates?

DIANA: $\mathcal{O} \left(\frac{\Delta_0 \left(1 + (1 + \omega) \sqrt{\omega/n} \right)}{\varepsilon^2} \right)$

4. MARINA

$$x^{k+1} = x^k - \gamma \cdot \frac{1}{n} \sum_{i=1}^n g_i^k$$

DIANA: $g_i^k = h_i^k + \mathcal{Q} \left(\nabla f_i(x^k) - h_i^k \right)$

$$h_i^{k+1} = h_i^k + \alpha \mathcal{Q} \left(\nabla f_i(x^k) - h_i^k \right)$$

$$x^{k+1} = x^k - \gamma \cdot \frac{1}{n} \sum_{i=1}^n g_i^k$$

DIANA: $g_i^k = h_i^k + \mathcal{Q} \left(\nabla f_i(x^k) - h_i^k \right)$

$$h_i^{k+1} = h_i^k + \alpha \mathcal{Q} \left(\nabla f_i(x^k) - h_i^k \right)$$

MARINA: $g_i^k = \begin{cases} \nabla f_i \left(x^k \right) & \text{w.p. } p \\ g^{k-1} + \mathcal{Q} \left(\nabla f_i \left(x^k \right) - \nabla f_i \left(x^{k-1} \right) \right) & \text{w.p. } 1-p \end{cases}$

$$x^{k+1} = x^k - \gamma \cdot \frac{1}{n} \sum_{i=1}^n g_i^k$$

DIANA: $g_i^k = h_i^k + \mathcal{Q} \left(\nabla f_i(x^k) - h_i^k \right)$

$$h_i^{k+1} = h_i^k + \alpha \mathcal{Q} \left(\nabla f_i(x^k) - h_i^k \right)$$

MARINA: $g_i^k = \begin{cases} \nabla f_i(x^k) & \text{w.p. } p \\ g^{k-1} + \mathcal{Q} \left(\nabla f_i(x^k) - \nabla f_i(x^{k-1}) \right) & \text{w.p. } 1-p \end{cases}$

typically small

p

$$x^{k+1} = x^k - \gamma \cdot \frac{1}{n} \sum_{i=1}^n g_i^k$$

DIANA: $g_i^k = h_i^k + \mathcal{Q} \left(\nabla f_i(x^k) - h_i^k \right)$

$$h_i^{k+1} = h_i^k + \alpha \mathcal{Q} \left(\nabla f_i(x^k) - h_i^k \right)$$

MARINA: $g_i^k = \begin{cases} \nabla f_i(x^k) \\ g^{k-1} + \mathcal{Q} \left(\nabla f_i(x^k) - \nabla f_i(x^{k-1}) \right) \end{cases}$

typically small

w.p. p

w.p. $1-p$

$$x^{k+1} = x^k - \gamma \cdot \frac{1}{n} \sum_{i=1}^n g_i^k$$

DIANA: $g_i^k = h_i^k + \mathcal{Q} \left(\nabla f_i(x^k) - h_i^k \right)$ vectors that devices have to send

$$h_i^{k+1} = h_i^k + \alpha \mathcal{Q} \left(\nabla f_i(x^k) - h_i^k \right)$$

MARINA: $g_i^k = \begin{cases} \nabla f_i(x^k) & \text{w.p. } p \\ g^{k-1} + \mathcal{Q} \left(\nabla f_i(x^k) - \nabla f_i(x^{k-1}) \right) & \text{w.p. } 1-p \end{cases}$

typically small

$$x^{k+1} = x^k - \gamma \cdot \frac{1}{n} \sum_{i=1}^n g_i^k = x^k - \gamma g^k$$

DIANA: $g_i^k = h_i^k + \mathcal{Q} \left(\nabla f_i(x^k) - h_i^k \right)$ vectors that devices have to send

$$h_i^{k+1} = h_i^k + \alpha \mathcal{Q} \left(\nabla f_i(x^k) - h_i^k \right)$$

MARINA: $g_i^k = \begin{cases} \nabla f_i(x^k) & \text{w.p. } p \\ g^{k-1} + \mathcal{Q} \left(\nabla f_i(x^k) - \nabla f_i(x^{k-1}) \right) & \text{w.p. } 1-p \end{cases}$

typically small

$$x^{k+1} = x^k - \gamma \cdot \frac{1}{n} \sum_{i=1}^n g_i^k = x^k - \gamma g^k$$

DIANA: $g_i^k = h_i^k + \mathcal{Q}(\nabla f_i(x^k) - h_i^k)$ vectors that devices have to send

$$h_i^{k+1} = h_i^k + \alpha \mathcal{Q}(\nabla f_i(x^k) - h_i^k) \quad \mathbb{E}[g^k | x^k] = \nabla f(x^k)$$

MARINA: $g_i^k = \begin{cases} \nabla f_i(x^k) & \text{w.p. } p \\ g^{k-1} + \mathcal{Q}(\nabla f_i(x^k) - \nabla f_i(x^{k-1})) & \text{w.p. } 1-p \end{cases}$

$$\mathbb{E}[g^k | x^k] \neq \nabla f(x^k)$$

Complexity Bound for MARINA

MARINA finds such \hat{x} that $\mathbb{E} \left[\|\nabla f(\hat{x})\|^2 \right] \leq \varepsilon^2$ after

Complexity Bound for MARINA

MARINA finds such \hat{x} that $\mathbb{E} \left[\|\nabla f(\hat{x})\|^2 \right] \leq \varepsilon^2$ after

$$\mathcal{O} \left(\frac{\Delta_0 (1 + \omega / \sqrt{n})}{\varepsilon^2} \right)$$

communication
rounds

Complexity Bound for MARINA

MARINA finds such \hat{x} that $\mathbb{E} \left[\|\nabla f(\hat{x})\|^2 \right] \leq \boxed{\varepsilon^2}$ after

Hides
numerical
factors and
smoothness
constants

$$\mathcal{O} \left(\frac{\Delta_0 (1 + \omega / \sqrt{n})}{\boxed{\varepsilon^2}} \right)$$

communication
rounds

Complexity Bound for MARINA

MARINA finds such \hat{x} that $\mathbb{E} \left[\|\nabla f(\hat{x})\|^2 \right] \leq \boxed{\varepsilon^2}$ after

Hides
numerical
factors and
smoothness
constants

$$\mathcal{O} \left(\frac{\Delta_0 \left(1 + \boxed{\omega / \sqrt{n}} \right)}{\boxed{\varepsilon^2}} \right)$$

communication
rounds

$$\boxed{\mathbb{E} \|Q(x) - x\|^2 \leq \omega \|x\|^2}$$

Complexity Bound for MARINA

MARINA finds such \hat{x} that $\mathbb{E} \left[\|\nabla f(\hat{x})\|^2 \right] \leq \boxed{\varepsilon^2}$ after

Hides
numerical
factors and
smoothness
constants

$$\mathcal{O} \left(\frac{\Delta_0 \left(1 + \boxed{\omega / \sqrt{n}} \right)}{\boxed{\varepsilon^2}} \right)$$

communication rounds

$$\mathbb{E} \|Q(x) - x\|^2 \leq \omega \|x\|^2$$

$$\Delta_0 = f(x^0) - f_*$$

Complexity Bound for MARINA

MARINA finds such \hat{x} that $\mathbb{E} \left[\|\nabla f(\hat{x})\|^2 \right] \leq \boxed{\varepsilon^2}$ after

Hides numerical factors and smoothness constants

$$\mathcal{O} \left(\frac{\Delta_0 \left(1 + \boxed{\omega / \sqrt{n}} \right)}{\varepsilon^2} \right)$$

communication rounds

$$\mathbb{E} \| \mathcal{Q}(x) - x \|^2 \leq \omega \| x \|^2$$

$$\Delta_0 = f(x^0) - f_*$$

$$p = \frac{1}{\omega + 1} = \Theta \left(\frac{\zeta_{\mathcal{Q}}}{d} \right)$$

Complexity Bound for MARINA

MARINA finds such \hat{x} that $\mathbb{E} \left[\|\nabla f(\hat{x})\|^2 \right] \leq \boxed{\varepsilon^2}$ after

Hides numerical factors and smoothness constants

$$\mathcal{O} \left(\frac{\Delta_0 \left(1 + \boxed{\omega / \sqrt{n}} \right)}{\varepsilon^2} \right)$$

communication rounds

$$\mathbb{E} \|Q(x) - x\|^2 \leq \omega \|x\|^2$$

$$\Delta_0 = f(x^0) - f_*$$

$$p = \frac{1}{\omega + 1} = \Theta \left(\frac{\zeta_Q}{d} \right)$$

assumption (holds for RandK, l2-quantization)

$$\zeta_Q = \sup_{x \in \mathbb{R}^d} \mathbb{E} [\|Q(x)\|_0]$$

expected density

Complexity Bounds for MARINA and DIANA

DIANA:

$$\mathcal{O} \left(\frac{\Delta_0 \left(1 + (1 + \omega) \sqrt{\omega/n} \right)}{\varepsilon^2} \right)$$

MARINA:

$$\mathcal{O} \left(\frac{\Delta_0 \left(1 + \omega / \sqrt{n} \right)}{\varepsilon^2} \right)$$

Complexity Bounds for MARINA and DIANA

DIANA:

$$\mathcal{O} \left(\frac{\Delta_0 \left(1 + (1 + \boxed{\omega}) \sqrt{\boxed{\omega}/n} \right)}{\varepsilon^2} \right)$$

MARINA:

$$\mathcal{O} \left(\frac{\Delta_0 \left(1 + \boxed{\omega}/\sqrt{n} \right)}{\varepsilon^2} \right)$$

5. MARINA and Variance Reduction

$$f_i(x) = \frac{1}{m} \sum_{j=1}^m f_{ij}(x)$$

$$f_i(x) = \frac{1}{m} \sum_{j=1}^m f_{ij}(x) \quad x^{k+1} = x^k - \gamma \cdot \frac{1}{n} \sum_{i=1}^n g_i^k$$

$$f_i(x) = \frac{1}{m} \sum_{j=1}^m f_{ij}(x) \quad x^{k+1} = x^k - \gamma \cdot \frac{1}{n} \sum_{i=1}^n g_i^k$$

MARINA: $g_i^k = \begin{cases} \nabla f_i(x^k) & \text{w.p. } p \\ g^{k-1} + \mathcal{Q}(\nabla f_i(x^k) - \nabla f_i(x^{k-1})) & \text{w.p. } 1-p \end{cases}$

$$f_i(x) = \frac{1}{m} \sum_{j=1}^m f_{ij}(x) \quad x^{k+1} = x^k - \gamma \cdot \frac{1}{n} \sum_{i=1}^n g_i^k$$

MARINA: $g_i^k = \begin{cases} \nabla f_i(x^k) & \text{w.p. } p \\ g^{k-1} + \mathcal{Q}(\nabla f_i(x^k) - \nabla f_i(x^{k-1})) & \text{w.p. } 1-p \end{cases}$

VR-MARINA: $g_i^k = \begin{cases} \nabla f_i(x^k) & \text{w.p. } p \\ g^{k-1} + \mathcal{Q}(\nabla f_{ij}(x^k) - \nabla f_{ij}(x^{k-1})) & \text{w.p. } 1-p \end{cases}$

$$f_i(x) = \frac{1}{m} \sum_{j=1}^m f_{ij}(x) \quad x^{k+1} = x^k - \gamma \cdot \frac{1}{n} \sum_{i=1}^n g_i^k$$

MARINA: $g_i^k = \begin{cases} \nabla f_i(x^k) & \text{w.p. } p \\ g^{k-1} + \mathcal{Q}(\nabla f_i(x^k) - \nabla f_i(x^{k-1})) & \text{w.p. } 1-p \end{cases}$

VR-MARINA: $g_i^k = \begin{cases} \nabla f_i(x^k) & \text{w.p. } p \\ g^{k-1} + \mathcal{Q}(\nabla f_{ij} \boxed{j}(x^k) - \nabla f_{ij} \boxed{j}(x^{k-1})) & \text{w.p. } 1-p \end{cases}$

$\boxed{j} \sim \{1, \dots, m\}$ uniformly at random

$$f_i(x) = \frac{1}{m} \sum_{j=1}^m f_{ij}(x) \quad x^{k+1} = x^k - \gamma \cdot \frac{1}{n} \sum_{i=1}^n g_i^k$$

MARINA: $g_i^k = \begin{cases} \nabla f_i(x^k) & \text{w.p. } p \\ g^{k-1} + \mathcal{Q}(\nabla f_i(x^k) - \nabla f_i(x^{k-1})) & \text{w.p. } 1-p \end{cases}$

VR-MARINA: $g_i^k = \begin{cases} \nabla f_i(x^k) & \text{w.p. } p \\ g^{k-1} + \mathcal{Q}(\nabla f_{ij} \boxed{j}(x^k) - \nabla f_{ij} \boxed{j}(x^{k-1})) & \text{w.p. } 1-p \end{cases}$

$\boxed{j} \sim \{1, \dots, \boxed{m}\}$ uniformly at random

Complexity Bound for VR-MARINA

VR-MARINA finds such \hat{x} that $\mathbb{E} \left[\|\nabla f(\hat{x})\|^2 \right] \leq \varepsilon^2$ after

Complexity Bound for VR-MARINA

VR-MARINA finds such \hat{x} that $\mathbb{E} \left[\|\nabla f(\hat{x})\|^2 \right] \leq \varepsilon^2$ after

$$\mathcal{O} \left(\frac{\Delta_0 \left(1 + \max\{\omega, \sqrt{(1+\omega)m}\} / \sqrt{n} \right)}{\varepsilon^2} \right)$$

communication
rounds/oracle
calls per node

Complexity Bound for VR-MARINA

VR-MARINA finds such \hat{x} that $\mathbb{E} \left[\|\nabla f(\hat{x})\|^2 \right] \leq \varepsilon^2$ after

Hides numerical factors and smoothness constants

$$\mathcal{O} \left(\frac{\Delta_0 \left(1 + \max\{\omega, \sqrt{(1+\omega)m}\} / \sqrt{n} \right)}{\varepsilon^2} \right)$$

communication rounds/oracle calls per node

Complexity Bound for VR-MARINA

VR-MARINA finds such \hat{x} that $\mathbb{E} \left[\|\nabla f(\hat{x})\|^2 \right] \leq \varepsilon^2$ after

Hides numerical factors and smoothness constants

$$\mathcal{O} \left(\frac{\Delta_0 \left(1 + \max\{\omega, \sqrt{(1+\omega)m}\} / \sqrt{n} \right)}{\varepsilon^2} \right)$$

communication rounds/oracle calls per node

Complexity Bound for VR-MARINA

VR-MARINA finds such \hat{x} that $\mathbb{E} \left[\|\nabla f(\hat{x})\|^2 \right] \leq \varepsilon^2$ after

Hides numerical factors and smoothness constants

$$\mathcal{O} \left(\frac{\Delta_0 \left(1 + \max\{\omega, \sqrt{(1+\omega)m}\}/\sqrt{n} \right)}{\varepsilon^2} \right)$$

communication rounds/oracle calls per node

$$\mathbb{E} \|Q(x) - x\|^2 \leq \omega \|x\|^2$$

Complexity Bound for VR-MARINA

VR-MARINA finds such \hat{x} that $\mathbb{E} \left[\|\nabla f(\hat{x})\|^2 \right] \leq \varepsilon^2$ after

Hides numerical factors and smoothness constants

$$\mathcal{O} \left(\frac{\Delta_0 \left(1 + \max\{\omega, \sqrt{(1+\omega)m}\}/\sqrt{n} \right)}{\varepsilon^2} \right)$$

communication rounds/oracle calls per node

$$\mathbb{E} \|Q(x) - x\|^2 \leq \omega \|x\|^2$$

$$\Delta_0 = f(x^0) - f_*$$

Complexity Bound for VR-MARINA

VR-MARINA finds such \hat{x} that $\mathbb{E} \left[\|\nabla f(\hat{x})\|^2 \right] \leq \varepsilon^2$ after

Hides numerical factors and smoothness constants

$$\mathcal{O} \left(\frac{\Delta_0 \left(1 + \max \{ \omega, \sqrt{(1 + \omega)m} \} / \sqrt{n} \right)}{\varepsilon^2} \right)$$

communication rounds/oracle calls per node

$$\mathbb{E} \|Q(x) - x\|^2 \leq \omega \|x\|^2$$

$$\Delta_0 = f(x^0) - f_*$$

$$p = \min \left\{ \frac{1}{\omega + 1}, \frac{1}{m + 1} \right\}$$

Complexity Bounds for MARINA and VR-MARINA

$$f_i(x) = \frac{1}{m} \sum_{j=1}^m f_{ij}(x)$$

communication rounds

oracle calls per node

MARINA:

$$\mathcal{O}\left(\frac{\Delta_0 (1 + \boxed{\omega}/\sqrt{n})}{\varepsilon^2}\right) \quad \mathcal{O}\left(\frac{\boxed{m}\Delta_0(1 + \boxed{\omega}/\sqrt{n})}{\varepsilon^2}\right)$$

VR-MARINA:

$$\mathcal{O}\left(\frac{\Delta_0 \left(1 + \max\{\boxed{\omega}, \sqrt{(1 + \boxed{\omega})\boxed{m}}\} / \sqrt{n}\right)}{\varepsilon^2}\right)$$

Complexity Bounds for VR-DIANA and VR-MARINA

$$f_i(x) = \frac{1}{m} \sum_{j=1}^m f_{ij}(x)$$

communication rounds

oracle calls per node

VR-DIANA:

$$\mathcal{O}\left(\frac{\Delta_0 \left(m^{2/3} + \omega\right) \sqrt{1 + \omega/n}}{\varepsilon^2}\right)$$

VR-MARINA:

$$\mathcal{O}\left(\frac{\Delta_0 \left(1 + \max\{\omega, \sqrt{(1 + \omega)m}/\sqrt{n}\}\right)}{\varepsilon^2}\right)$$

6. MARINA and Partial Participation

PP-MARINA

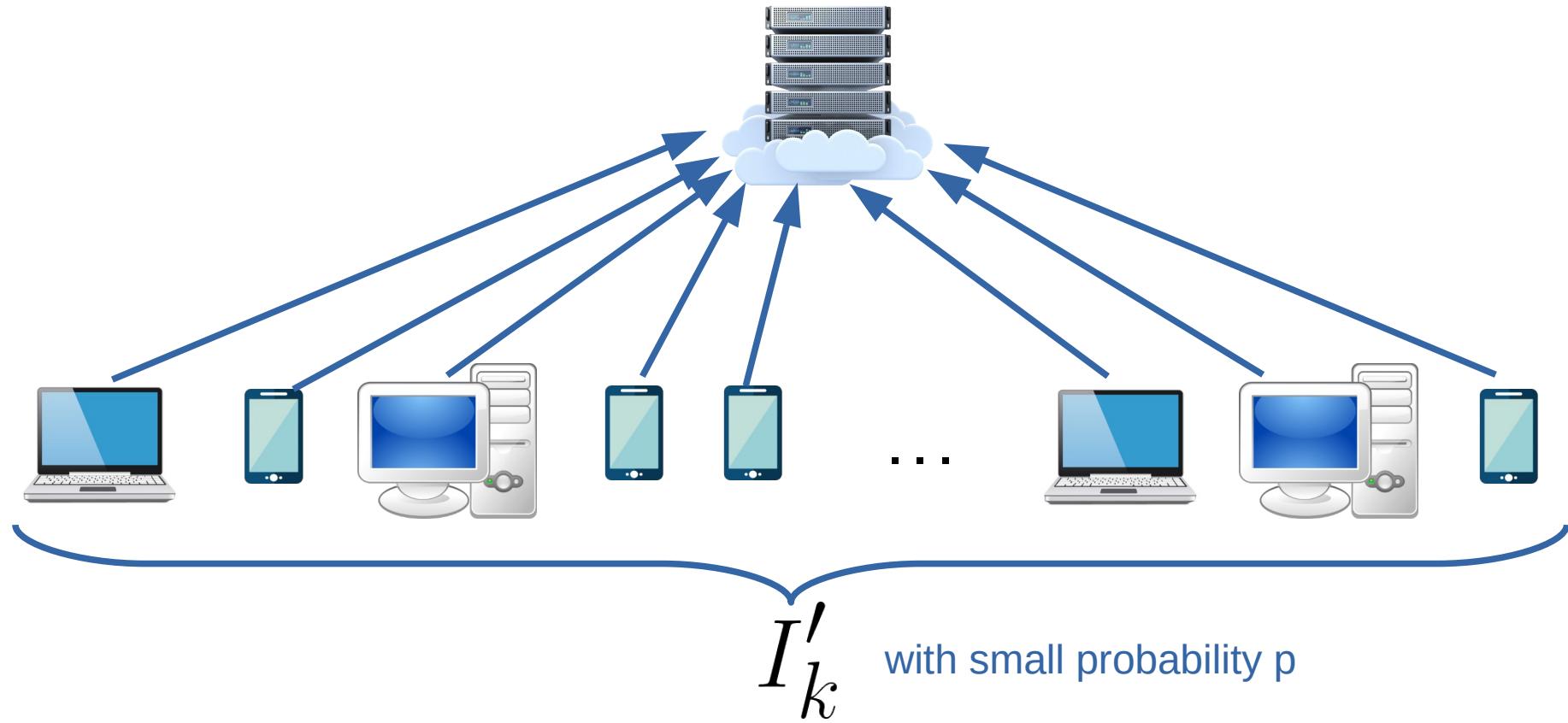


PP-MARINA

 I'_k

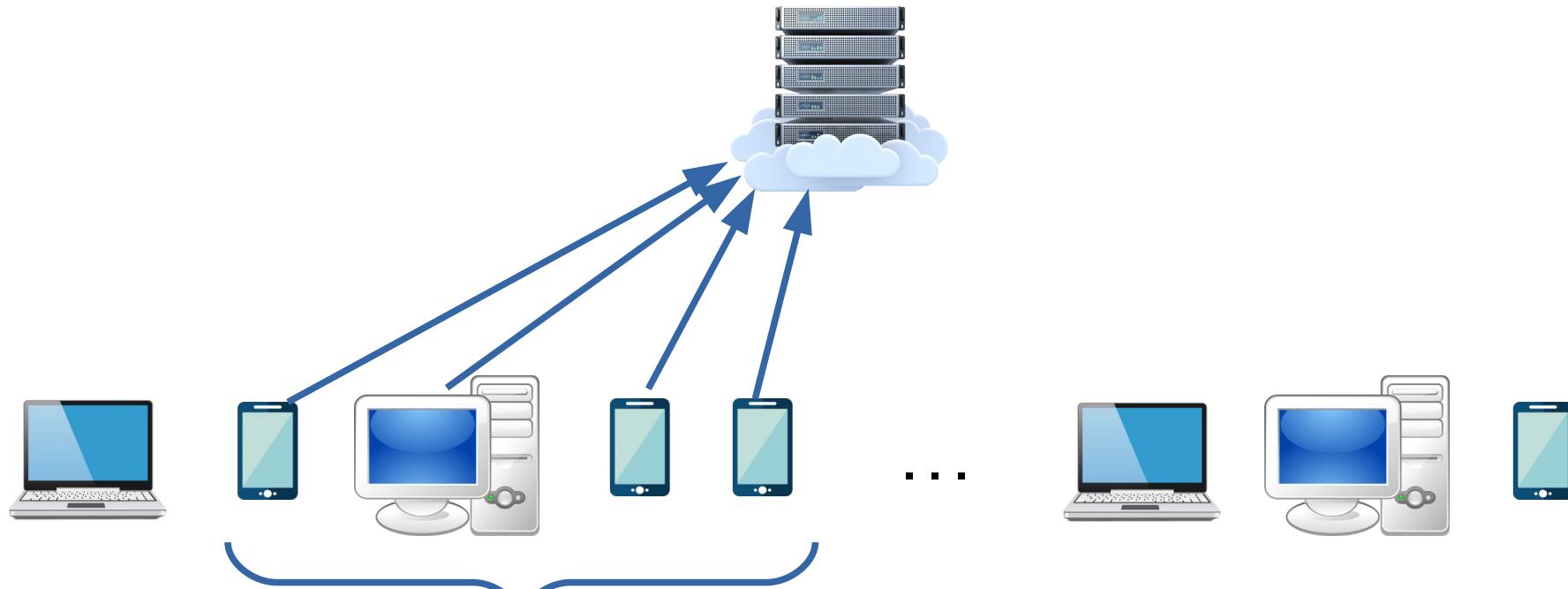
- the set of clients participating in communication on step k

PP-MARINA



I'_k - the set of clients participating in communication on step k

PP-MARINA



I'_k with large probability $1 - p$

I'_k - the set of clients participating in communication on step k

PP-MARINA

$$x^{k+1} = x^k - \gamma g^k$$

PP-MARINA

$$x^{k+1} = x^k - \gamma g^k$$

$$g^k = \begin{cases} \nabla f(x^k) & \text{if } c_k = 1 \\ g^{k-1} + \frac{1}{r} \sum_{i_k \in I'_k} \mathcal{Q}(\nabla f_{i_k}(x^k) - \nabla f_{i_k}(x^{k-1})) & \text{if } c_k = 0 \end{cases}$$

PP-MARINA

$$x^{k+1} = x^k - \gamma g^k$$

$$c_k \sim \text{Be}(p)$$

$$g^k = \begin{cases} \nabla f\left(x^k\right) & \text { if } c_k=1 \\ g^{k-1}+\frac{1}{r} \sum_{i_k \in I'_k} \mathcal{Q}\left(\nabla f_{i_k}\left(x^k\right)-\nabla f_{i_k}\left(x^{k-1}\right)\right) & \text { if } c_k=0 \end{cases}$$

PP-MARINA

$$x^{k+1} = x^k - \gamma g^k$$

full participation & uncompressed communication

$$c_k \sim \text{Be}(p)$$

$$g^k = \begin{cases} \boxed{\nabla f(x^k)} & \text{if } c_k = 1 \\ g^{k-1} + \frac{1}{r} \sum_{i_k \in I'_k} \mathcal{Q} \left(\nabla f_{i_k}(x^k) - \nabla f_{i_k}(x^{k-1}) \right) & \text{if } c_k = 0 \end{cases}$$

PP-MARINA

$$x^{k+1} = x^k - \gamma g^k$$

full participation & uncompressed communication

$$c_k \sim \text{Be}(p)$$

$$g^k = \begin{cases} \nabla f(x^k) \\ g^{k-1} + \frac{1}{r} \sum_{i_k \in I'_k} \mathcal{Q}(\nabla f_{i_k}(x^k) - \nabla f_{i_k}(x^{k-1})) \end{cases}$$

if $c_k = 1$

if $c_k = 0$

partial participation & compressed communication

$$|I'_k| = r$$

Complexity Bound for PP-MARINA

PP-MARINA finds such \hat{x} that $\mathbb{E} \left[\|\nabla f(\hat{x})\|^2 \right] \leq \varepsilon^2$ after

Complexity Bound for PP-MARINA

PP-MARINA finds such \hat{x} that $\mathbb{E} \left[\|\nabla f(\hat{x})\|^2 \right] \leq \varepsilon^2$ after

$$\mathcal{O} \left(\frac{\Delta_0 (1 + (1 + \omega) \sqrt{n}/r)}{\varepsilon^2} \right)$$

communication
rounds

Complexity Bound for PP-MARINA

PP-MARINA finds such \hat{x} that $\mathbb{E} \left[\|\nabla f(\hat{x})\|^2 \right] \leq \boxed{\varepsilon^2}$ after

Hides
numerical
factors and
smoothness
constants

$$\rightarrow \mathcal{O} \left(\frac{\Delta_0 (1 + (1 + \omega) \sqrt{n/r})}{\boxed{\varepsilon^2}} \right)$$

communication rounds

Complexity Bound for PP-MARINA

PP-MARINA finds such \hat{x} that $\mathbb{E} \left[\|\nabla f(\hat{x})\|^2 \right] \leq \boxed{\varepsilon^2}$ after

Hides numerical factors and smoothness constants

$$\mathcal{O} \left(\frac{\Delta_0 (1 + (1 + \boxed{\omega}) \sqrt{n/r})}{\boxed{\varepsilon^2}} \right)$$

communication rounds

$$\mathbb{E} \|Q(x) - x\|^2 \leq \omega \|x\|^2$$

Complexity Bound for PP-MARINA

PP-MARINA finds such \hat{x} that $\mathbb{E} \left[\|\nabla f(\hat{x})\|^2 \right] \leq \boxed{\varepsilon^2}$ after

Hides numerical factors and smoothness constants $\rightarrow \mathcal{O} \left(\frac{\Delta_0 (1 + (1 + \boxed{\omega}) \sqrt{n/r})}{\boxed{\varepsilon^2}} \right)$ communication rounds

$$\mathbb{E} \|Q(x) - x\|^2 \leq \omega \|x\|^2$$

$$\Delta_0 = f(x^0) - f_*$$

Complexity Bound for PP-MARINA

PP-MARINA finds such \hat{x} that $\mathbb{E} \left[\|\nabla f(\hat{x})\|^2 \right] \leq \boxed{\varepsilon^2}$ after

Hides numerical factors and smoothness constants $\rightarrow \mathcal{O} \left(\frac{\Delta_0 (1 + (1 + \boxed{\omega}) \boxed{\sqrt{n/r}})}{\varepsilon^2} \right)$ communication rounds

$$\mathbb{E} \|Q(x) - x\|^2 \leq \omega \|x\|^2$$

$$\Delta_0 = f(x^0) - f_*$$

$$p = \frac{\boxed{r}}{n(\omega + 1)} = \Theta \left(\frac{\boxed{r} \zeta_Q}{nd} \right)$$

Complexity Bound for PP-MARINA

PP-MARINA finds such \hat{x} that $\mathbb{E} \left[\|\nabla f(\hat{x})\|^2 \right] \leq \boxed{\varepsilon^2}$ after

Hides numerical factors and smoothness constants $\rightarrow \mathcal{O} \left(\frac{\Delta_0 (1 + (1 + \boxed{\omega}) \boxed{\sqrt{n/r}})}{\varepsilon^2} \right)$ communication rounds

$$\mathcal{O} \left(\frac{\Delta_0 (1 + (1 + \boxed{\omega}) \boxed{\sqrt{n/r}})}{\varepsilon^2} \right)$$

$$\mathbb{E} \|Q(x) - x\|^2 \leq \omega \|x\|^2$$

$$\Delta_0 = f(x^0) - f_*$$

$$p = \frac{r}{n(\omega + 1)} = \Theta \left(\frac{r \zeta_Q}{nd} \right)$$

\uparrow assumption (holds for RandK, l2-quantization)

\uparrow expected density

$$\zeta_Q = \sup_{x \in \mathbb{R}^d} \mathbb{E} [\|Q(x)\|_0]$$

Complexity Bounds for PP-MARINA and MARINA

MARINA: $\mathcal{O} \left(\frac{\Delta_0 (1 + \omega / \sqrt{n})}{\varepsilon^2} \right)$

PP-MARINA: $\mathcal{O} \left(\frac{\Delta_0 (1 + (1 + \omega) \sqrt{n/r})}{\varepsilon^2} \right)$

7. Experiments

The Problem

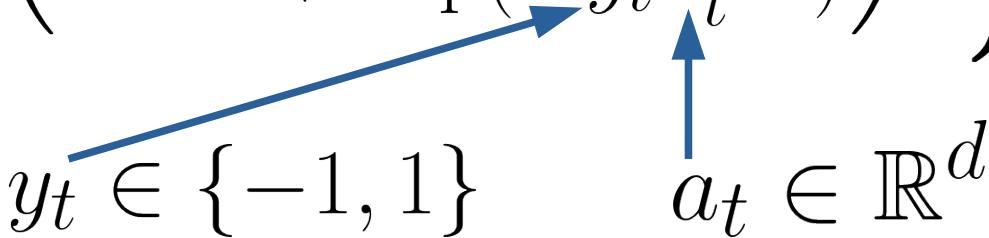
$$\min_{x \in \mathbb{R}^d} \left\{ f(x) = \frac{1}{N} \sum_{t=1}^N \left(1 - \frac{1}{1 + \exp(-y_t a_t^\top x)} \right)^2 \right\}$$

The Problem

$$\min_{x \in \mathbb{R}^d} \left\{ f(x) = \frac{1}{N} \sum_{t=1}^N \left(1 - \frac{1}{1 + \exp(-y_t a_t^\top x)} \right)^2 \right\}$$

A diagram illustrating the components of the loss function. Two blue arrows point upwards from the labels $y_t \in \{-1, 1\}$ and $a_t \in \mathbb{R}^d$ to the terms y_t and $a_t^\top x$ in the equation, respectively.

The Problem

$$\min_{x \in \mathbb{R}^d} \left\{ f(x) = \frac{1}{N} \sum_{t=1}^N \left(1 - \frac{1}{1 + \exp(-y_t a_t^\top x)} \right)^2 \right\}$$


- The dataset was split into 5 equal parts among 5 clients

The Problem

$$\min_{x \in \mathbb{R}^d} \left\{ f(x) = \frac{1}{N} \sum_{t=1}^N \left(1 - \frac{1}{1 + \exp(-y_t a_t^\top x)} \right)^2 \right\}$$

$y_t \in \{-1, 1\}$ $a_t \in \mathbb{R}^d$

- The dataset was split into 5 equal parts among 5 clients
- Theoretical stepsizes

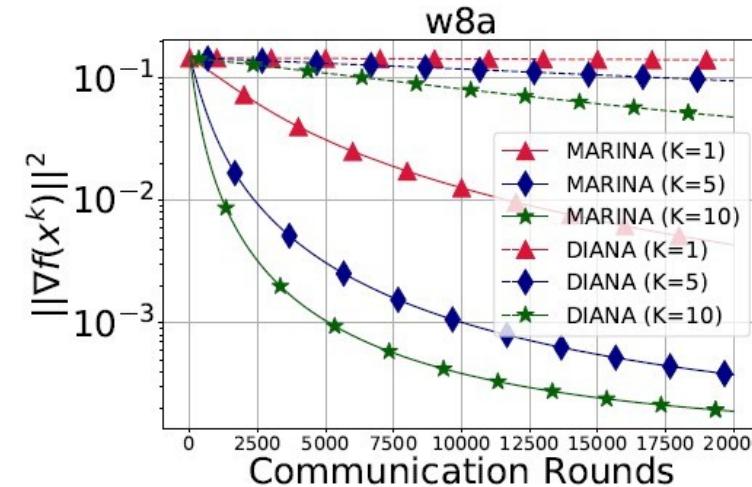
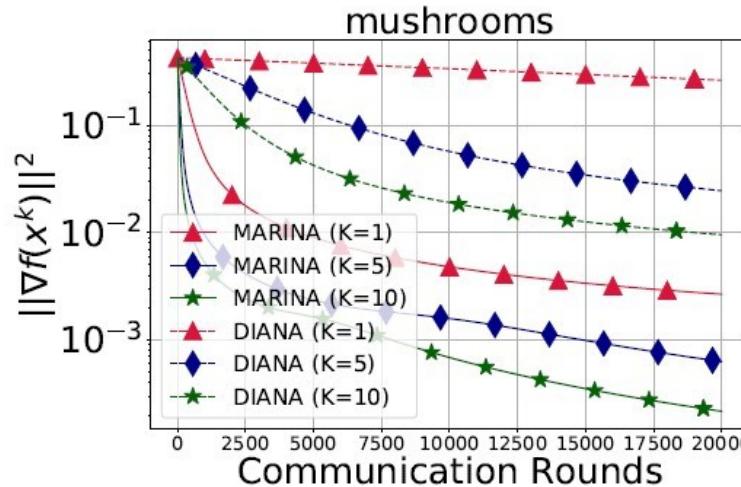
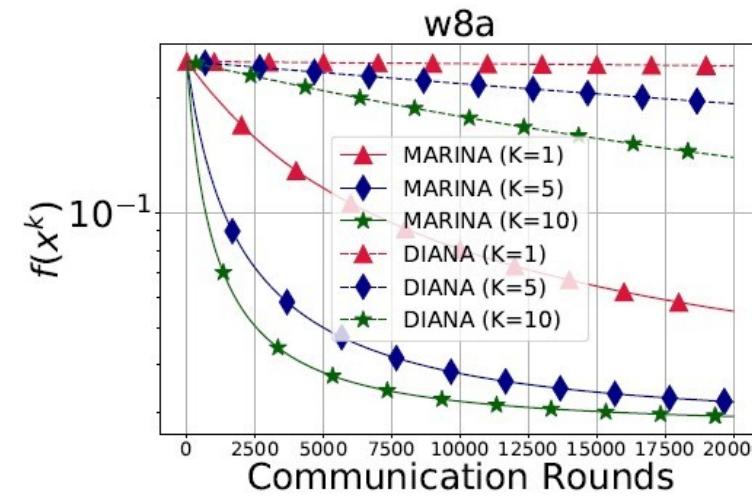
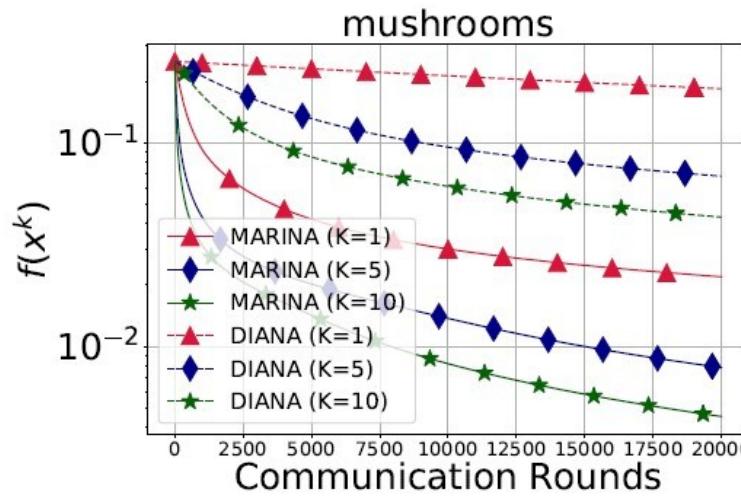
The Problem

$$\min_{x \in \mathbb{R}^d} \left\{ f(x) = \frac{1}{N} \sum_{t=1}^N \left(1 - \frac{1}{1 + \exp(-y_t a_t^\top x)} \right)^2 \right\}$$

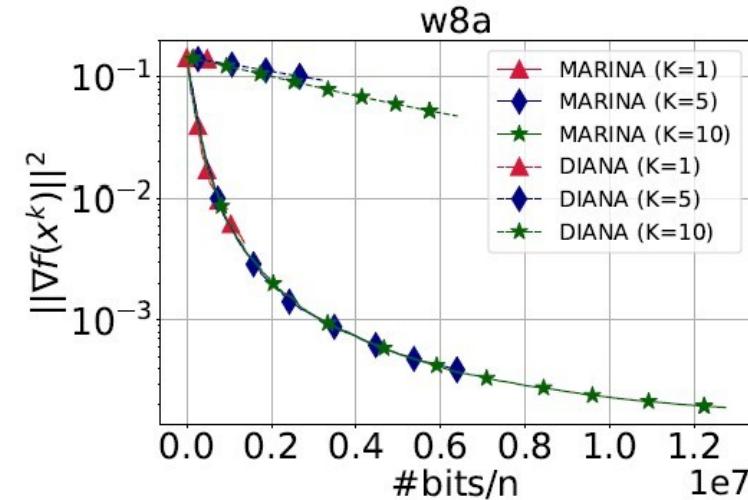
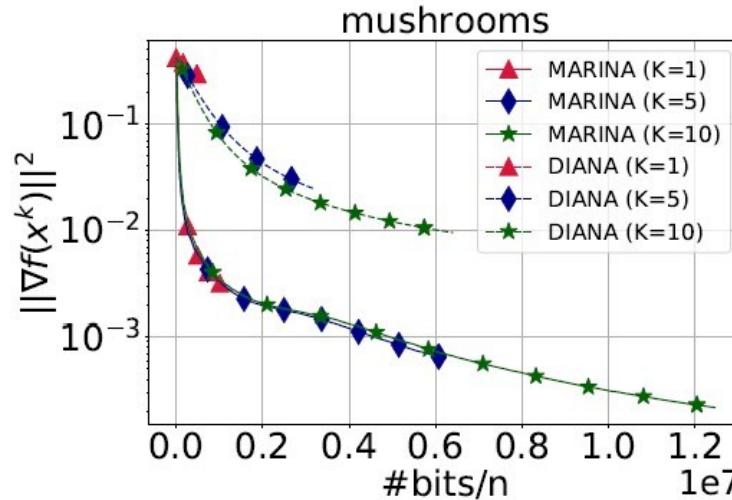
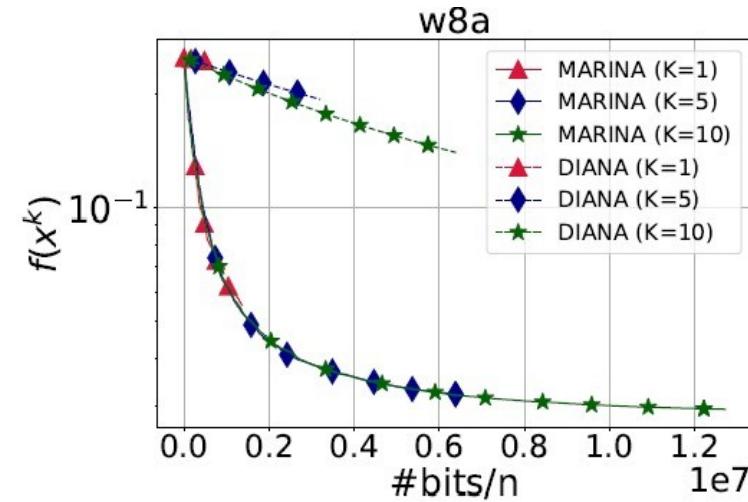
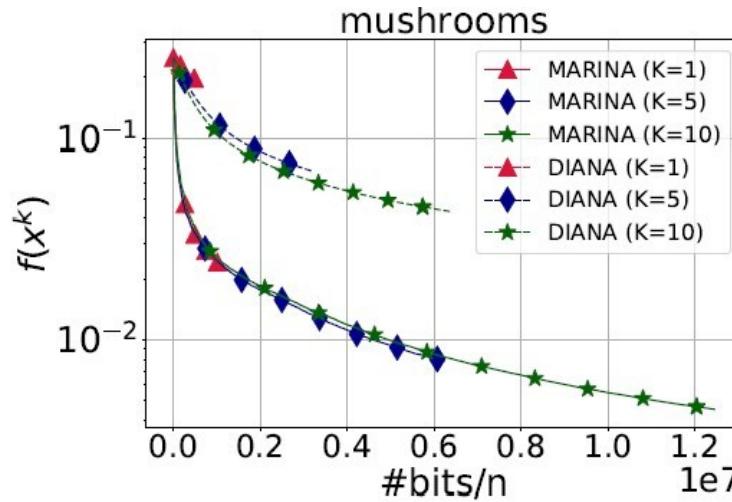
$y_t \in \{-1, 1\}$ $a_t \in \mathbb{R}^d$

- The dataset was split into 5 equal parts among 5 clients
- Theoretical stepsizes
- For stochastic methods batchsize was 1% of local data

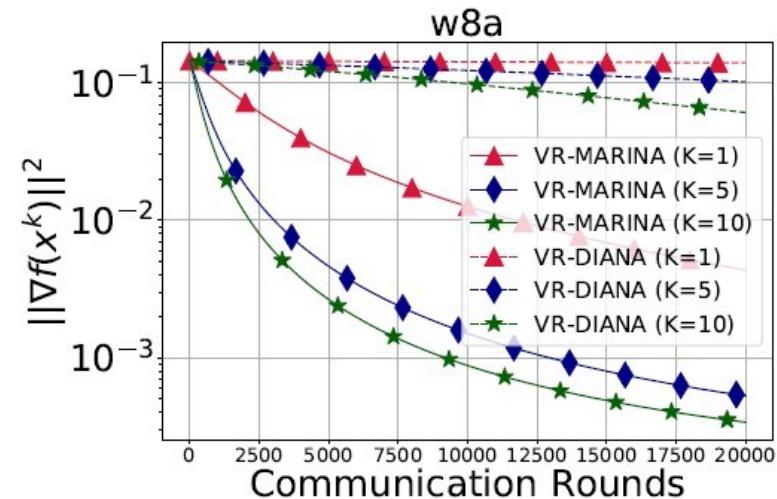
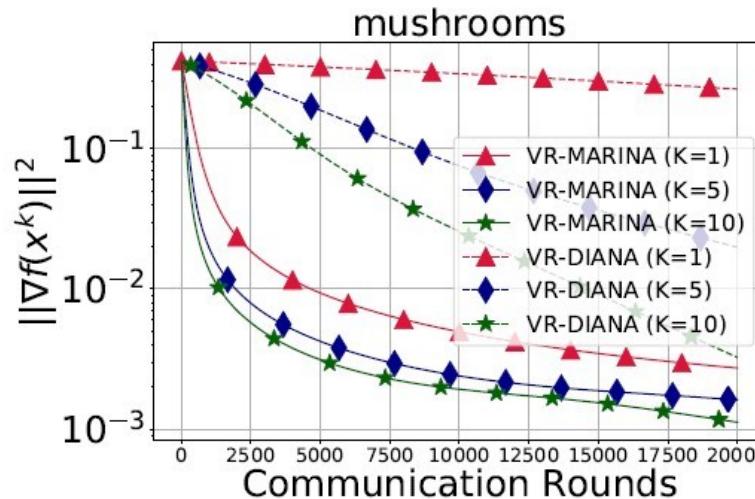
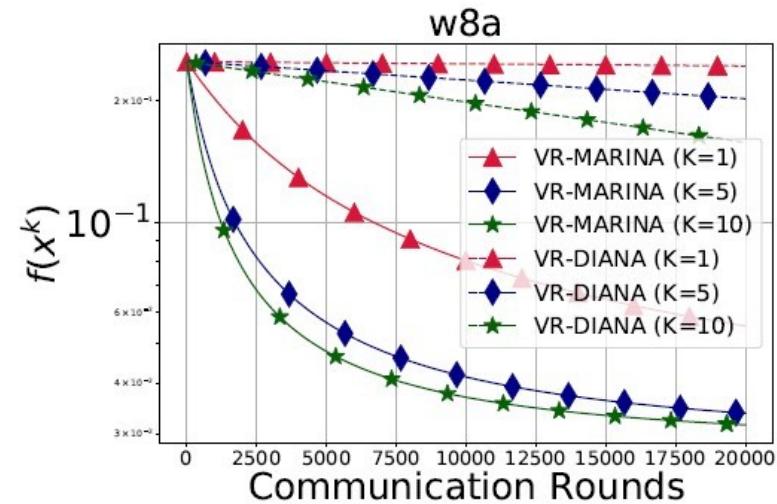
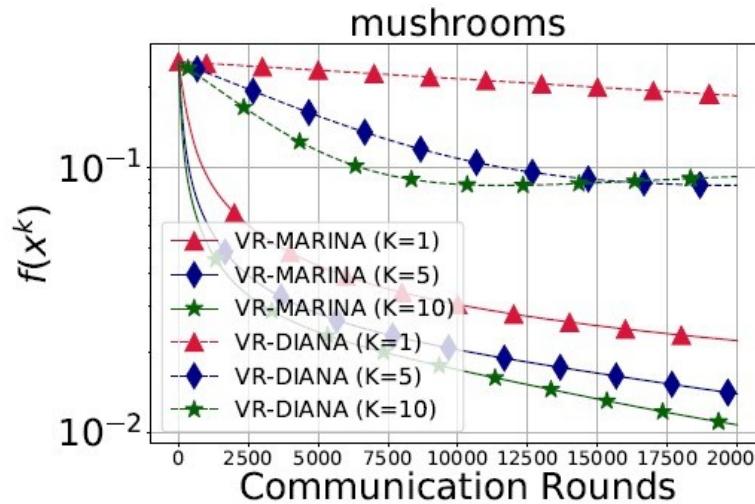
MARINA vs DIANA (RandK)



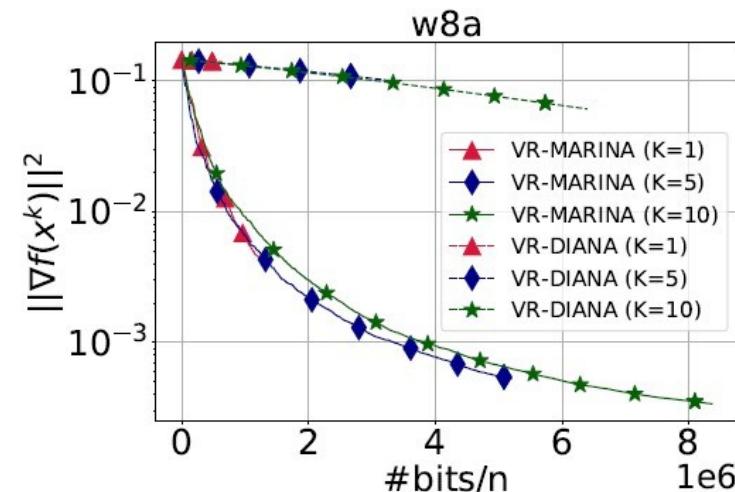
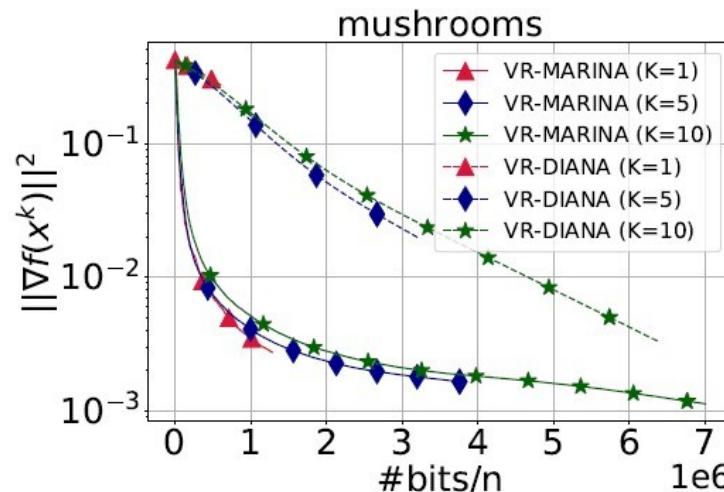
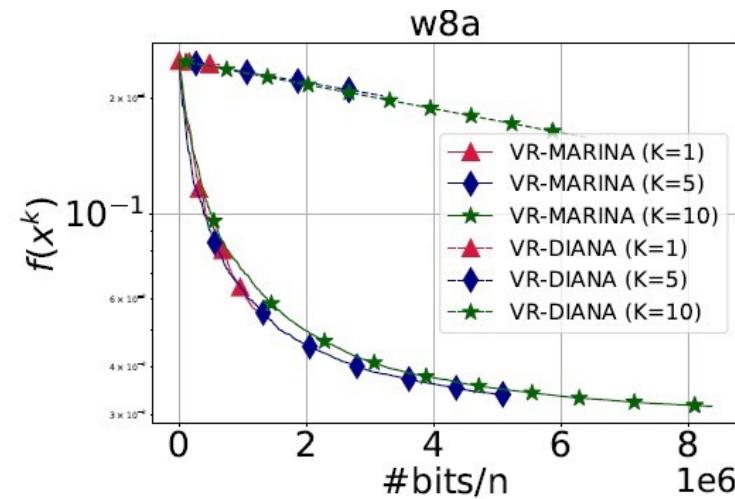
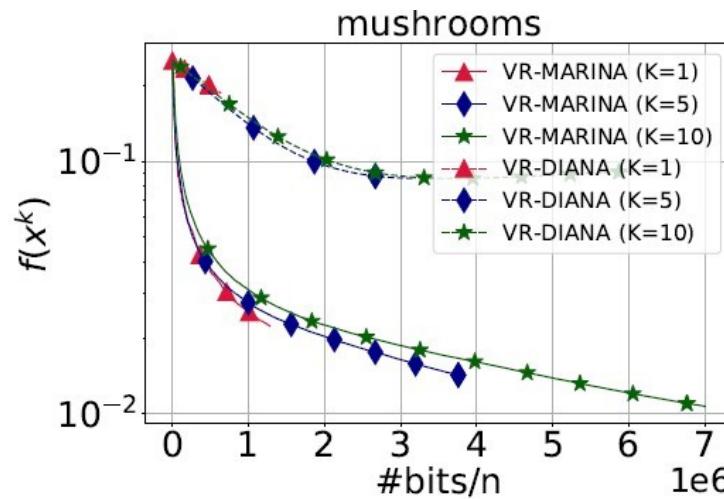
MARINA vs DIANA (RandK)



VR-MARINA vs VR-DIANA (RandK)



VR-MARINA vs VR-DIANA (RandK)



8. Extra Results

In the paper, we also have:

- Mini-batched version of Variance Reduced MARINA

In the paper, we also have:

- Mini-batched version of Variance Reduced MARINA
- Variance Reduced MARINA for the problems with $f_i(x) = \mathbf{E}_{\xi_i \sim \mathcal{D}_i} [f_{\xi_i}(x)]$

In the paper, we also have:

- Mini-batched version of Variance Reduced MARINA
- Variance Reduced MARINA for the problems with $f_i(x) = \mathbf{E}_{\xi_i \sim \mathcal{D}_i} [f_{\xi_i}(x)]$
- Rates under Polyak- Lojasiewicz Condition

In the paper, we also have:

- Mini-batched version of Variance Reduced MARINA
- Variance Reduced MARINA for the problems with $f_i(x) = \mathbf{E}_{\xi_i \sim \mathcal{D}_i} [f_{\xi_i}(x)]$
- Rates under Polyak- Lojasiewicz Condition
- Explicit dependencies on smoothness constants, non-uniform sampling

In the paper, we also have:

- Mini-batched version of Variance Reduced MARINA
- Variance Reduced MARINA for the problems with $f_i(x) = \mathbf{E}_{\xi_i \sim \mathcal{D}_i} [f_{\xi_i}(x)]$
- Rates under Polyak- Lojasiewicz Condition
- Explicit dependencies on smoothness constants, non-uniform sampling
- Simple proofs