

## 1. Problem

We are solving the distributed optimization problem:

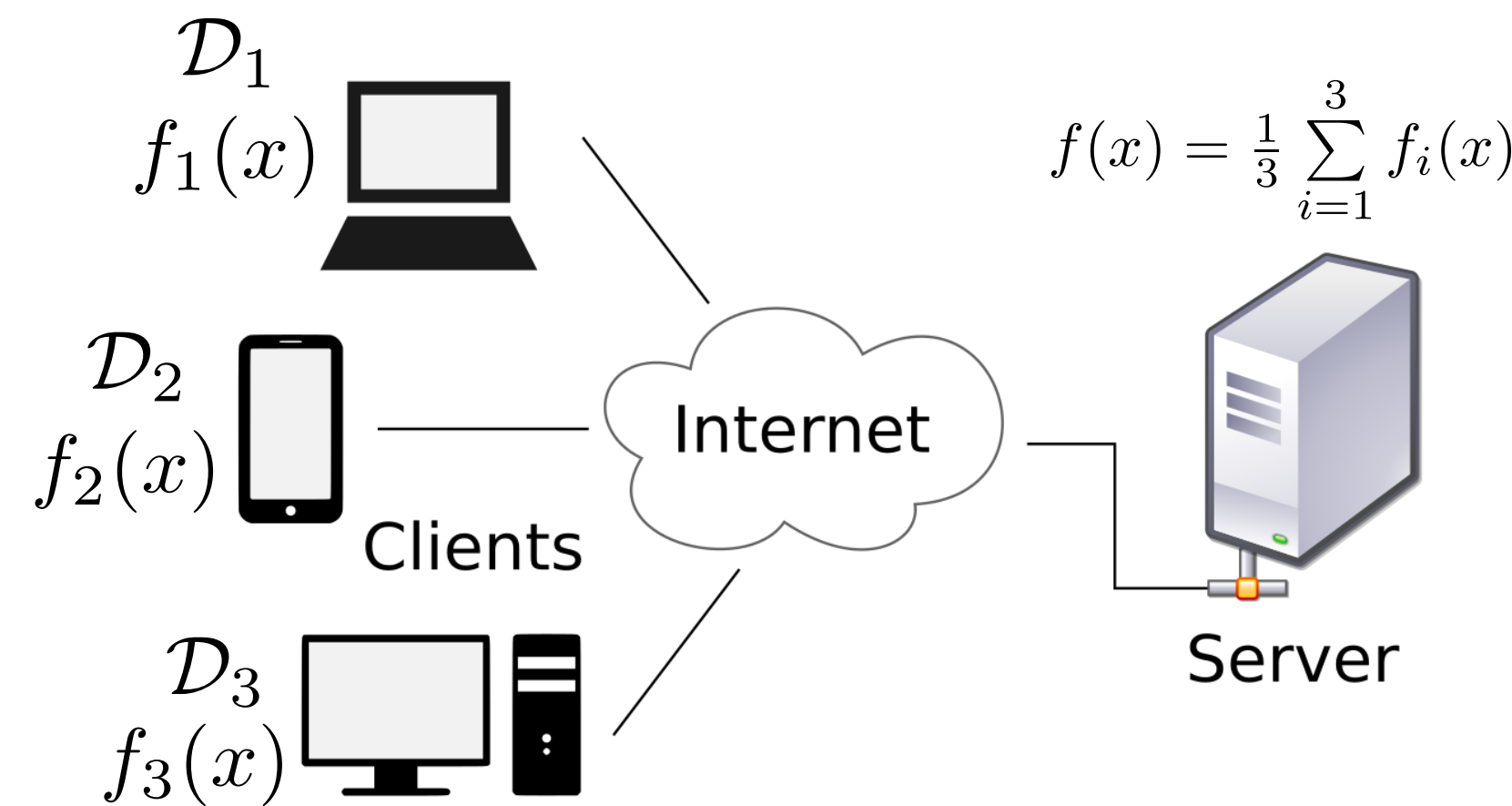
$$\min_{x \in \mathbb{R}^d} f(x) \stackrel{\text{def}}{=} \frac{1}{n} \sum_{i=1}^n f_i(x)$$

# devices / machines

# model parameters / features

The functions  $f_i$  can be arbitrarily heterogeneous/different

## 2. Technical Setup



Communication is the bottleneck in both directions!

## 3. Main Baseline: Vanilla GD Method

$$x^{t+1} = x^t - \gamma \nabla f(x^t) = x^t - \gamma \frac{1}{n} \sum_{i=1}^n \nabla f_i(x^t)$$

Recap: Communication Complexity of GD:

$$d \times \mathcal{O} \left( \frac{L}{\mu} \log \frac{1}{\varepsilon} \right)$$

# of coordinates # of GD iterations

Recap: Convergence of GD

GD returns an  $\mathcal{E}$ -solution (i.e.,  $\|x^T - x^*\|^2 \leq \varepsilon$ ) after

$$T = \frac{2L}{\mu} \log \left( \frac{\|x^0 - x^*\|^2}{\varepsilon} \right) = \mathcal{O} \left( \frac{L}{\mu} \log \frac{1}{\varepsilon} \right) \text{ iterations.}$$

## 4. New SOTA: EF21-P + DIANA and EF21-P + DCGD

### Algorithm 1 EF21-P + DIANA

- 1: **Parameters:** learning rates  $\gamma > 0$  (for learning the model) and  $\beta > 0$  (for learning the gradient shifts); initial model  $x^0 \in \mathbb{R}^d$  (stored on the server and the workers); initial gradient shifts  $h_1^0, \dots, h_n^0 \in \mathbb{R}^d$  (stored on the workers); average of the initial gradient shifts  $h^0 = \frac{1}{n} \sum_{i=1}^n h_i^0$  (stored on the server); initial model shift  $w^0 = x^0 \in \mathbb{R}^d$  (stored on the server and the workers)
- 2: **for**  $t = 0, 1, \dots, T-1$  **do**
- 3:   **for**  $i = 1, \dots, n$  **in parallel do**
- 4:      $m_i^t = \mathcal{C}^D(\nabla f_i(w^t) - h_i^t)$
- 5:     Send compressed message  $m_i^t$  to the server
- 6:      $h_i^{t+1} = h_i^t + \beta m_i^t$
- 7:   **end for**
- 8:    $m^t = \frac{1}{n} \sum_{i=1}^n m_i^t$
- 9:    $h^{t+1} = h^t + \beta m^t$
- 10:    $g^t = h^t + m^t$
- 11:    $x^{t+1} = x^t - \gamma g^t$
- 12:    $p^{t+1} = \mathcal{C}^P(x^{t+1} - w^t)$
- 13:    $w^{t+1} = w^t + p^{t+1}$
- 14:   Broadcast compressed message  $p^{t+1}$  to all  $n$  workers
- 15:   **for**  $i = 1, \dots, n$  **in parallel do**
- 16:      $w^{t+1} = w^t + p^{t+1}$
- 17:   **end for**
- 18: **end for**

DIANA (idea from [1])

EF21-P (new idea)

## 5. Contributions

1. **EF21-P + DIANA** provides **new state-of-the-art convergence rate** for distributed optimization in the **strongly convex** and **general convex** regimes.
2. **EF21-P + DIANA** is the **first method** supporting bidirectional compression whose server-to-workers and workers-to-server communication complexity is **no worse (can be much better!)** than that of vanilla **GD** in the **strongly convex** and **general convex** regimes:

### Communication Complexity of EF21-P + DIANA

$$K \times \mathcal{O} \left( \left( \frac{L}{\alpha \mu} + \omega + \frac{\omega L_{\max}}{n \mu} \right) \log \frac{1}{\varepsilon} \right)$$

# of coordinates # of EF21-P + DIANA iterations

where  $\omega$  and  $\alpha$  are compression parameters.

3. In the **nonconvex regime**, **EF21-P + DCGD** provides the **new state-of-the-art convergence rate** in the **low accuracy regimes** ( $\mathcal{E}$  is small or the # of workers  $n$  is large). We provide examples of optimization problems where **EF21-P + DCGD** achieves **new state-of-the-art convergence rate** even in the **high accuracy regime**.

## References

[1] Konstantin Mishchenko, Eduard Gorbunov, Martin Takáč, Peter Richtárik  
Distributed learning with compressed gradient differences  
arXiv:1901.09269

## 6. Main Tools: Compression Operators

### Unbiased compressor

$$\mathbb{E}[\mathcal{C}(x)] = x, \quad \mathbb{E}[\|\mathcal{C}(x) - x\|^2] \leq \omega \|x\|^2$$

**Example: RandK**     $K = 1$

$$\omega = \frac{d}{K} - 1 = \frac{4}{1} - 1 = 3$$

$$\begin{bmatrix} -3 \\ 2 \\ 6 \\ -8 \end{bmatrix} \xrightarrow{c} \underbrace{\begin{bmatrix} 4 \\ 0 \\ 0 \\ 0 \end{bmatrix}}_{d/K} \times \begin{bmatrix} 0 \\ 2 \\ 0 \\ 0 \end{bmatrix}$$

Random entry

### Biased compressor

$$\mathbb{E}[\|\mathcal{C}(x) - x\|^2] \leq (1 - \alpha) \|x\|^2 \quad \alpha \in (0, 1]$$

**Example: TopK**     $K = 1$

$$\alpha = \frac{K}{d} = \frac{1}{4}$$

$$\begin{bmatrix} -3 \\ 2 \\ 6 \\ -8 \end{bmatrix} \xrightarrow{c} \begin{bmatrix} 0 \\ 0 \\ 0 \\ -8 \end{bmatrix}$$

Largest entry (in magnitude)

## 7. New SOTA in the Convex Setting

For method comparison, we use the notation  $\omega_w \equiv \omega$ ,  $\omega_s \equiv 1/\alpha - 1$

Method	# Communication Rounds
EF (No server-to-worker compression!) (Seide et al., 2014)	$\Omega \left( (1 + \omega_w) \frac{L_{\max}}{\mu} \right)$
DIANA (No server-to-worker compression!) (Mishchenko et al., 2019)	$(1 + \frac{\omega_w}{n}) \frac{L_{\max}}{\mu} + \omega_w$
Dore, Artemis (Liu et al., 2020) (Philippenko & Dieuleveut, 2020)	$\Omega \left( \frac{\omega_s \omega_w}{n} \frac{L_{\max}}{\mu} \right)$
MCM (Philippenko & Dieuleveut, 2021)	$\Omega \left( \left( \omega_s^{3/2} + \frac{\omega_s \omega_w^{1/2}}{\sqrt{n}} + \frac{\omega_w}{n} \right) \frac{L_{\max}}{\mu} \right)$
EF21-P + DIANA (new)	$(1 + \omega_s) \frac{L}{\mu} + \frac{\omega_w}{n} \frac{L_{\max}}{\mu} + \omega_w$

The relationship between  $\omega_s$  and  $\omega_w$  in EF21-P + DIANA is **linear!**

## 8. New SOTA in the Non-Convex Setting

MCM (Philippenko & Dieuleveut, 2021)	$\frac{\omega_s^{3/2}}{\varepsilon} + \frac{\omega_s \omega_w^{1/2}}{\sqrt{n\varepsilon}} + \frac{\omega_w}{n\varepsilon}$
CD-Adam (Wang et al., 2022)	$\frac{\sqrt{d} \max\{\omega_s, \omega_w\}^4}{\varepsilon^2}$
EF21-BC (Fatkhullin et al., 2021)	$\frac{\omega_w \omega_s}{\varepsilon}$
EF21-P + DCGD (new)	$\frac{\omega_w}{n\varepsilon^2} + \frac{\omega_s}{\varepsilon}$
EF21-P + DCGD (new)	$\frac{D\omega_w}{n\varepsilon} + \frac{\omega_s}{\varepsilon}$ (strong-growth assumption)