

# A Better Alternative to Error-Feedback for Communication-Efficient Distributed Learning

Samuel Horváth and Peter Richtárik

King Abdullah University of Science and Technology

## Problem setup

We consider **distributed optimization problems** of the form

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

- $n$  is the number of nodes,
- $f_i: \mathbb{R}^d \rightarrow \mathbb{R}$  is a smooth differentiable loss function composed of data stored on worker  $i$ .

## Communication Bottleneck

In distributed training, **model updates (or gradient vectors) have to be exchanged in each iteration**. Due to the size of the communicated messages for commonly considered deep models (Alistarh et al., 2016), this represents **significant bottleneck** of the whole optimization procedure. To reduce the amount of data that has to be transmitted, **communication compression** is one of the popular approaches. Considering both practice and theory, compression operators can be split into two groups: **biased** and **unbiased**.

### Definition ("Unbiased")

A randomized mapping  $C: \mathbb{R}^d \rightarrow \mathbb{R}^d$  is an **unbiased compression operator (unbiased compressor)** if there exists  $\beta \in [0, 1]$  such that

$$\mathbb{E}[C(x)] = x, \quad \mathbb{E}\|C(x) - x\|^2 \leq \beta \|x\|^2, \quad x \in \mathbb{R}^d.$$

If this holds, we will for simplicity write  $C \in \mathcal{U}(\beta)$ .

### Definition ("Biased")

A (possibly) randomized mapping  $C: \mathbb{R}^d \rightarrow \mathbb{R}^d$  is a **general compression operator (general compressor)** if there exists  $\beta > 0$  and  $\gamma \in [0, 1]$  such that

$$\mathbb{E}^h \|C(x) - x\|^2 \leq \beta \|x\|^2, \quad x \in \mathbb{R}^d.$$

If this holds, we will for simplicity write  $C \in \mathcal{C}(\beta, \gamma)$ .

## Induced Compressor vs. Error-Feedback (EF) [Example]

### Assumptions:

Biased and Unbiased Compressors:  $C_1 \in \mathcal{C}(\beta)$  and  $C_2 \in \mathcal{U}(\beta)$ .  
 $f$  is over-parametrized (Vaswani et al., 2019) and  $\mu$ -quasi convex, i.e.

$$f(x) = f(x^*) + \frac{\mu}{2} \|x - x^*\|^2, \quad x \in \mathbb{R}^d.$$

where  $x^*$  is the optimal solution of  $f$  and  $f = \frac{1}{n} \sum_{i=1}^n f_i$ ,  $f_i$ 's are  $L$ -smooth.

### Construction [on worker $i$ ]:

#### Induced Compressor

obtain  $g_i^k [E^h g_i^k = \nabla f_i(x^k)]$

$$h_i^k = C_1^k(g_i^k)$$

send  $h_i^k$  to master

[no need to keep track of errors]

### Convergence Rates $[E^h f(\bar{x}^T) - f^*]$ :

$$O\left(\frac{1}{n} + 1\right) L \exp\left(-\frac{\mu T}{4nL}\right)$$

#### Error-Feedback

obtain  $g_i^k [E^h g_i^k = \nabla f_i(x^k)]$

$$h_i^k = C_1^k(g_i^k + e_i^k)$$

send  $h_i^k$  to master

$$e_i^{k+1} = g_i^k + e_i^k - h_i^k$$

$$O\left(\frac{1}{n} + 1\right) L \exp\left(-\frac{\mu T}{4nL}\right)$$

## Contributions

- Induced Compressor.** When used with EF framework, biased compressors (e.g., Top-K) can often achieve superior performance when compared to their unbiased counterparts (e.g., Rand-K), which is attributed to their low empirical variance. Our key contribution is the development of a simple but remarkably effective alternative (described above), which we argue leads to better and more versatile methods both in theory and practice.
- Better Theory for DCSGD.** We provide a new and tighter theoretical analysis of DCSGD under weaker assumptions.
- Partial Participation.** We extend our results to obtain the first convergence guarantee for partial participation with arbitrary distributions over nodes.

## Motivational Example

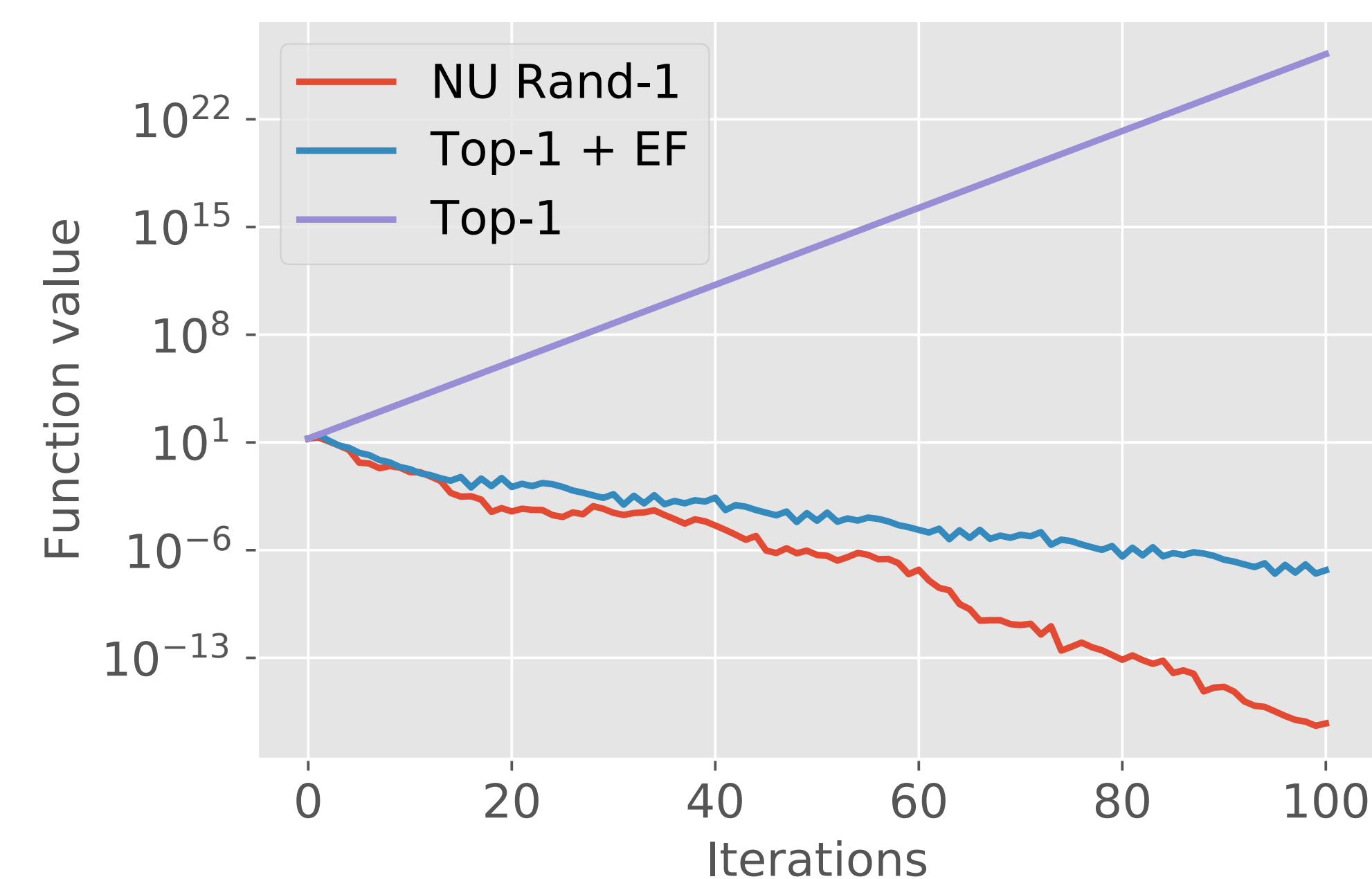


Figure: Comparison of Top-1 (+ EF) and NU Rand-1 on Example 1 from Beznosikov et al., 2020.

## Numerical Experiments

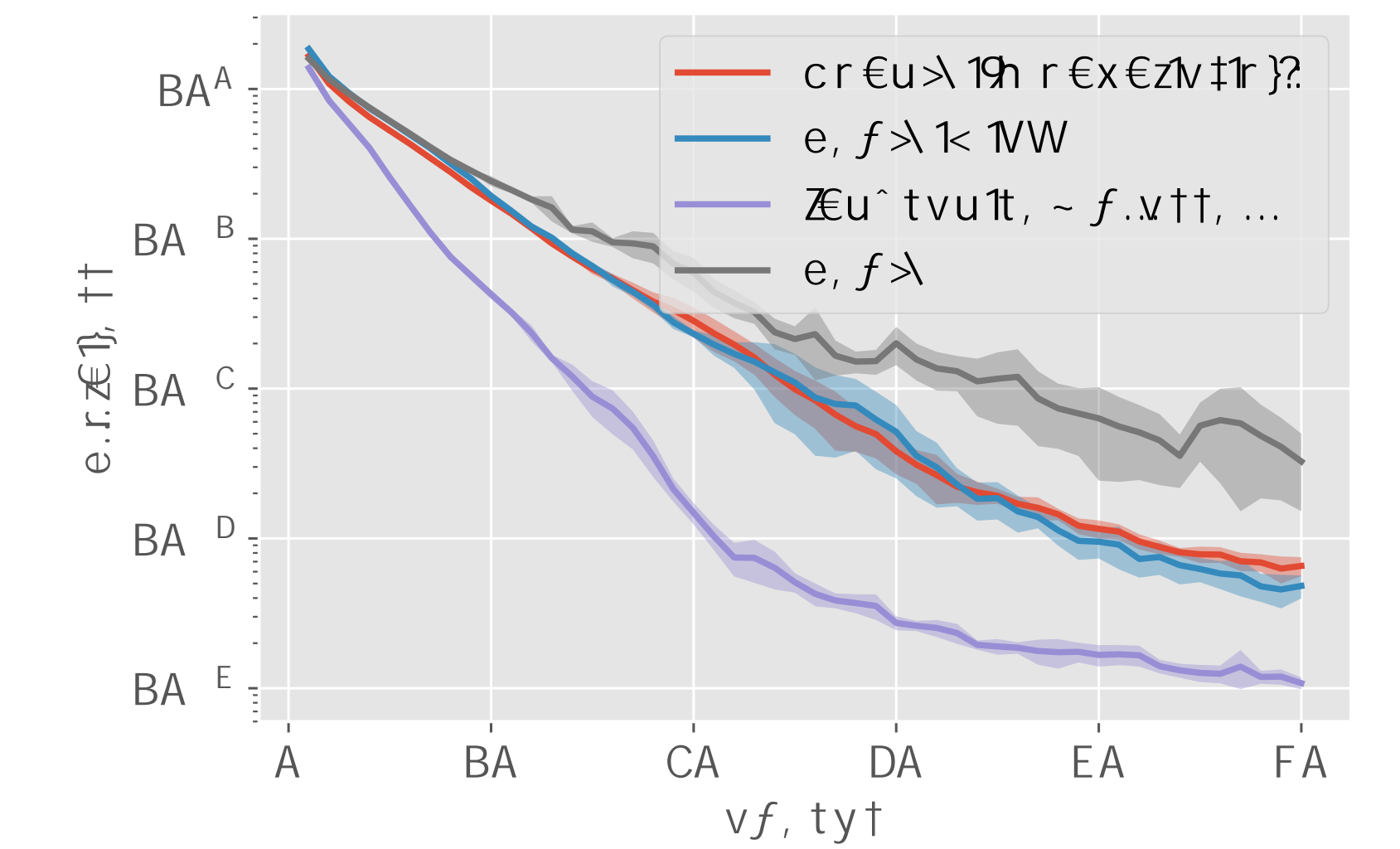
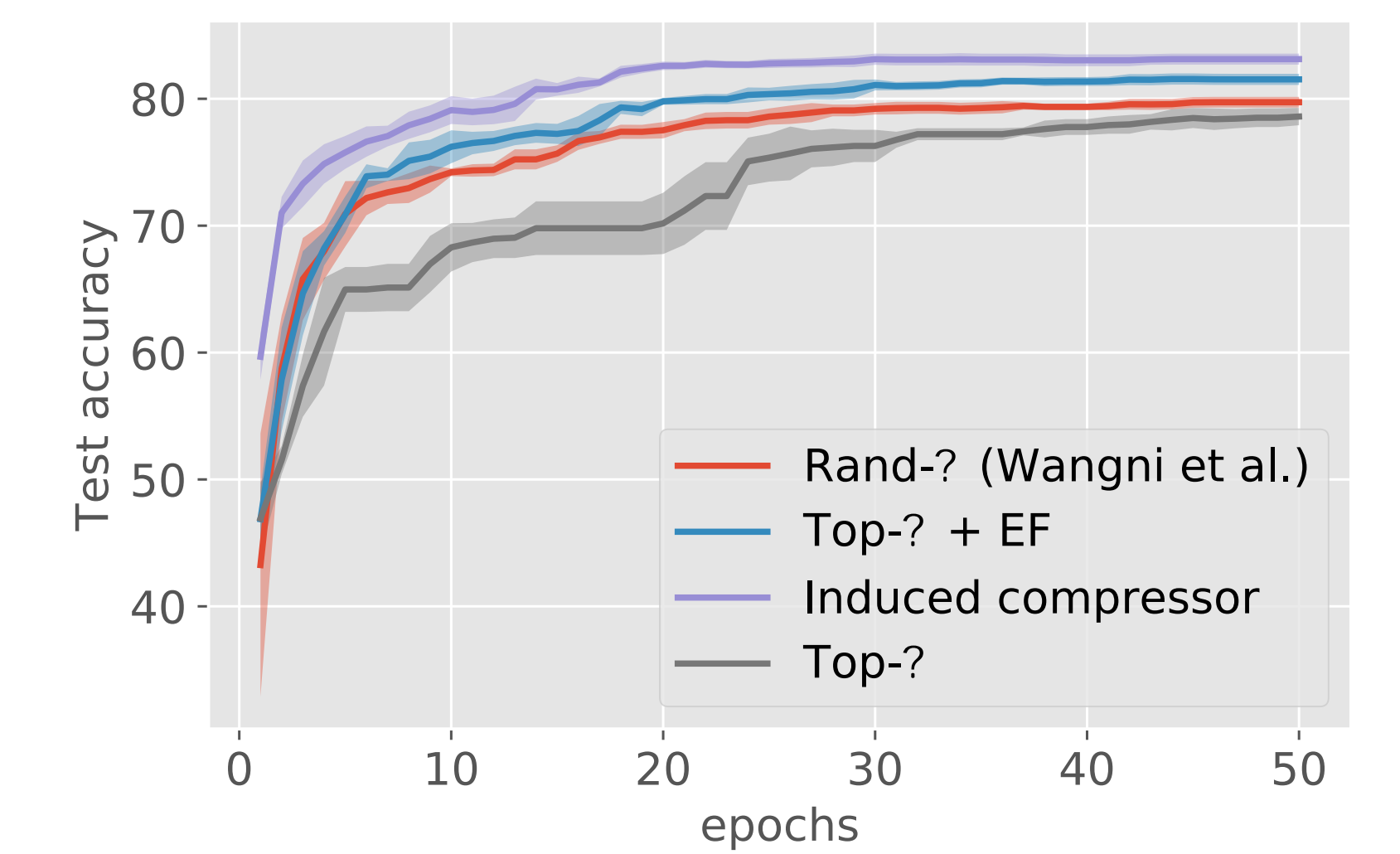


Figure: Comparison of different sparsification techniques on CIFAR10 with Resnet18.

## Contact Information

- Paper: <https://arxiv.org/pdf/2002.05359.pdf>
- Email: samuel.horvath@kaust.edu.sa

