

Gradient Compressed Sensing: A Query-Efficient Gradient Estimator for High-Dimensional Zeroth-Order Optimization

Ruizhong Qiu[†] Hanghang Tong[†]

{rq5, htong}@illinois.edu

[†] UIUC ILLINOIS DEAC



ICML
International Conference
On Machine Learning

HIGHLIGHTS

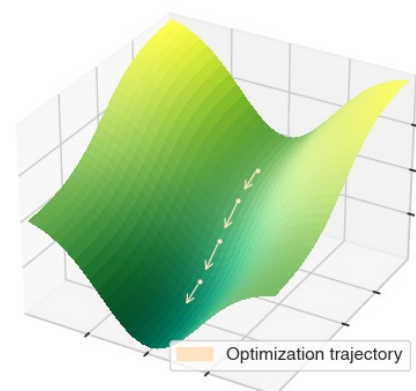
- **Query-efficient gradient estimator: GraCe**
 - Only $O(s \log \log(d/s))$ queries per step
 - Still $O(1/T)$ convergence for nonconvex optimization
- **Relaxed sparsity assumption**
 - We assume **approximately** s -sparse gradients
 - Weaker than *exact sparsity* and *compressibility*
- **Improvement of Indyk–Price–Woodruff (IPW)**
 - We reduce its constant by a factor of nearly **4300**
 - Via our *dependent random partition* technique
- **Strong empirical performance**
 - Evaluated under **10000**-dimensional functions
 - Significantly outperforms 12 existing methods

BACKGROUND

➤ High-dim zeroth-order optimization (ZOO)

$$\min_{x \in \mathbb{R}^d} f(x)$$

- High-dimensional space \mathbb{R}^d
- Possibly nonconvex function $f: \mathbb{R}^d \rightarrow \mathbb{R}$
- Only queries $f(x)$; no gradients $\nabla f(x)$
- Aim to use as few queries as possible



➤ Our main assumption: ρ -approximate s -sparsity

$$\max_{I \subseteq [d]: |I|=s} \|\nabla_I f(x)\|_2^2 \geq \rho \|\nabla f(x)\|_2^2, \quad \forall x$$

- Approximation ratio $0 < \rho \leq 1$; sparsity $1 \leq s \leq d$

➤ Other standard assumptions:

- Lower boundedness: $f_* := \inf_x f(x) > -\infty$
- Lipschitz continuity: $|f(x + u) - f(x)| \leq L_0 \|u\|_2$
- Lipschitz smoothness: $\|\nabla f(x + u) - \nabla f(x)\|_2 \leq L_1 \|u\|_2$

PROPOSED METHOD: GraCe

- **Core idea:** Locate large-gradient dimensions
 - Employ the **IPW** algorithm from compressed sensing
 - **Adaptively encode** dimension information into queries
- **Motivating case:** Approx. 1-sparse gradient
 - Suppose that some $j \in [d]$ has sufficiently large $\frac{|\nabla_j f(x)|}{\|\nabla_{[d] \setminus j} f(x)\|_2}$
 - Let $u'_i := \epsilon$ and $v'_i := \epsilon \cdot i$. Then $\frac{f(x+v') - f(x)}{f(x+u') - f(x)} \approx \frac{\epsilon \cdot j \cdot \nabla_j f(x)}{\epsilon \cdot \nabla_j f(x)} = j$
 - For instance, $f(x) := 0.1x_1 + x_2$ gives $\frac{f(0+v') - f(0)}{f(0+u') - f(0)} = \frac{2.1\epsilon}{1.1\epsilon} \approx 2$
- **General case: IPW + dependent random partition**
 - Divide dimensions into $O(s)$ groups S of fixed size $O(d/s)$
 - Each group has only 1 large-gradient dimension j w.h.p.
 - Use $O(\log \log |S|)$ **adaptive** queries to locate $j \in S$

THEORETICAL GUARANTEES

- **Query complexity:** $O(s \log \log(d/s))$
 - Given any $x \in \mathbb{R}^d$, any $\epsilon > 0$, and any $0 < \alpha < \rho$
 - $O(s \log \log(d/s))$ queries can find a gradient estimate g s.t.

$$\|g\|_2 \leq \|\nabla f(x)\|_2 + O(\epsilon)$$

$$\mathbb{E}[\langle \nabla f(x), g \rangle | x] \geq \alpha \|\nabla f(x)\|_2^2 - O(\epsilon)$$
- **Rate of convergence:** $O(1/T)$ for nonconvex ZOO
 - Given any initial point $x_1 \in \mathbb{R}^d$, any step size $0 < \eta < \rho/L_1$, any $0 < \beta < 1$, and any $\Delta > 0$, under suitable hyperparams
 - Suppose gradient descent + GraCe yields points x_2, x_3, \dots
 - With probability $\geq 1 - \beta$, for all $T \geq 1$ simultaneously,

$$\min_{t=1, \dots, T} \|\nabla f(x_t)\|_2^2 \leq \frac{1 + \frac{2(1-L_1\eta)}{L_1\eta\beta}}{\eta - \frac{L_1\eta^2}{2}} (f(x_1) - f_*) + \Delta$$

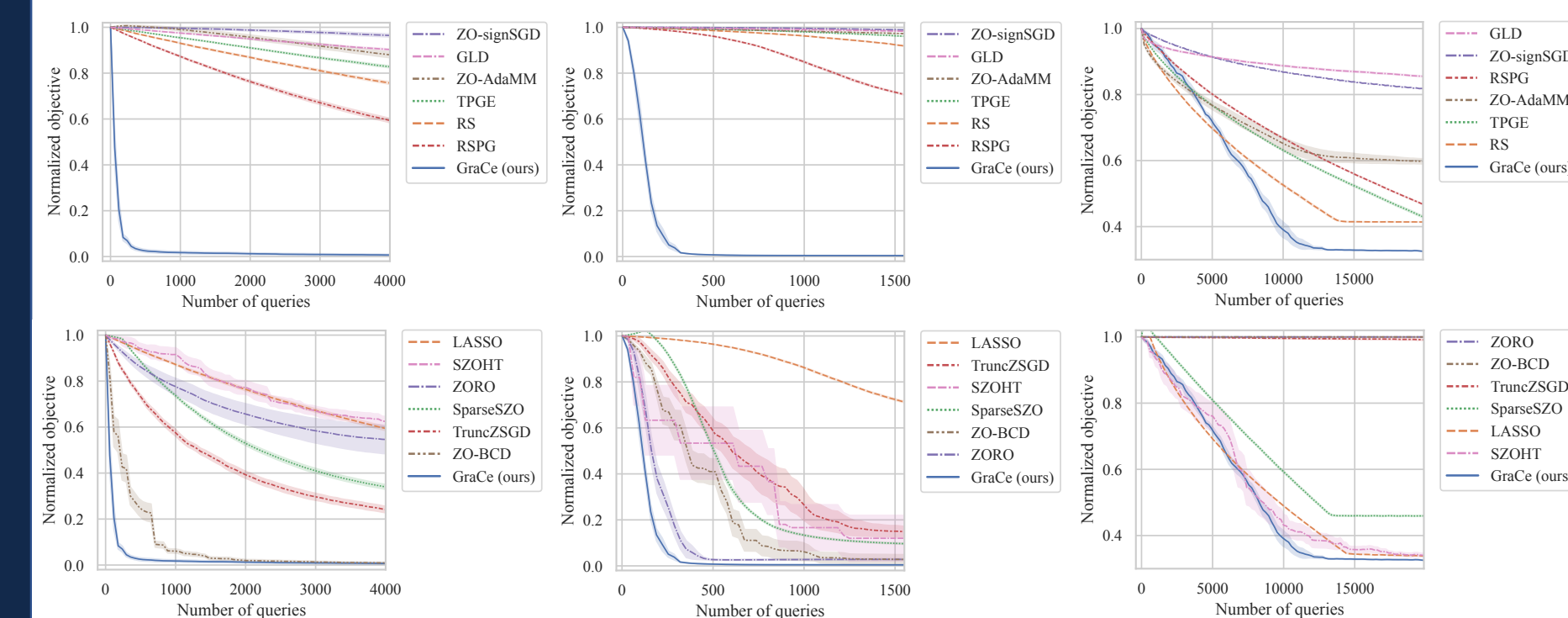
COMPARISON w/ BASELINES

➤ Theoretical comparison:

Type	Method	Queries per step	Rate of convergence
Full Gradient	RS (Ghadimi & Lan, 2012)	$O(1)$	$\mathbb{E}[\ \nabla f(x_\tau)\ _2^2] \leq O(\frac{\sqrt{d}}{\sqrt{T}} + \frac{d}{T})$
	TPGE (Duchi et al., 2015)	$O(1)$	$\mathbb{E}[\ \nabla f(x_\tau)\ _2^2] \leq O(\frac{\sqrt{d}}{\sqrt{T}})$
	RSPG (Ghadimi et al., 2016)	$O(q)$	$\mathbb{E}[\ \nabla f(x_\tau)\ _2^2] \leq O(\frac{d}{q} + \frac{d^2}{qT})$
	ZO-signSGD (Liu et al., 2019)	$O(bq)$	$\mathbb{E}[\ \nabla f(x_\tau)\ _2] \leq O(\frac{\sqrt{d}\sqrt{q+d}}{\sqrt{bq}} + \frac{\sqrt{d}}{\sqrt{T}})$
	ZO-AdaMM (Chen et al., 2019)	$O(1)$	$\mathbb{E}[\ \nabla f(x_\tau)\ _2^2] \leq O(\frac{d}{T} + \frac{d^2}{T})$
Sparse Gradient	ZORO (Cai et al., 2022)	$O(s \log \frac{d}{s})$	$\ \nabla f(x_\tau)\ _2^2 \leq O(\frac{1}{T})$ w.h.p.
	GraCe (ours)	$O(s \log \log \frac{d}{s})$	$\ \nabla f(x_\tau)\ _2^2 \leq O(\frac{1}{T})$ w.h.p.

➤ Empirical comparison:

- Evaluated under **10000**-dimensional functions
- 2 synthetic functions (left) and 1 real-world function (right)



➤ Conclusions:

- GraCe consistently outperforms 12 existing ZOO methods
- GraCe achieves **fastest** convergence with **fewest** queries

ACKNOWLEDGEMENTS



C3.ai Digital Transformation Institute