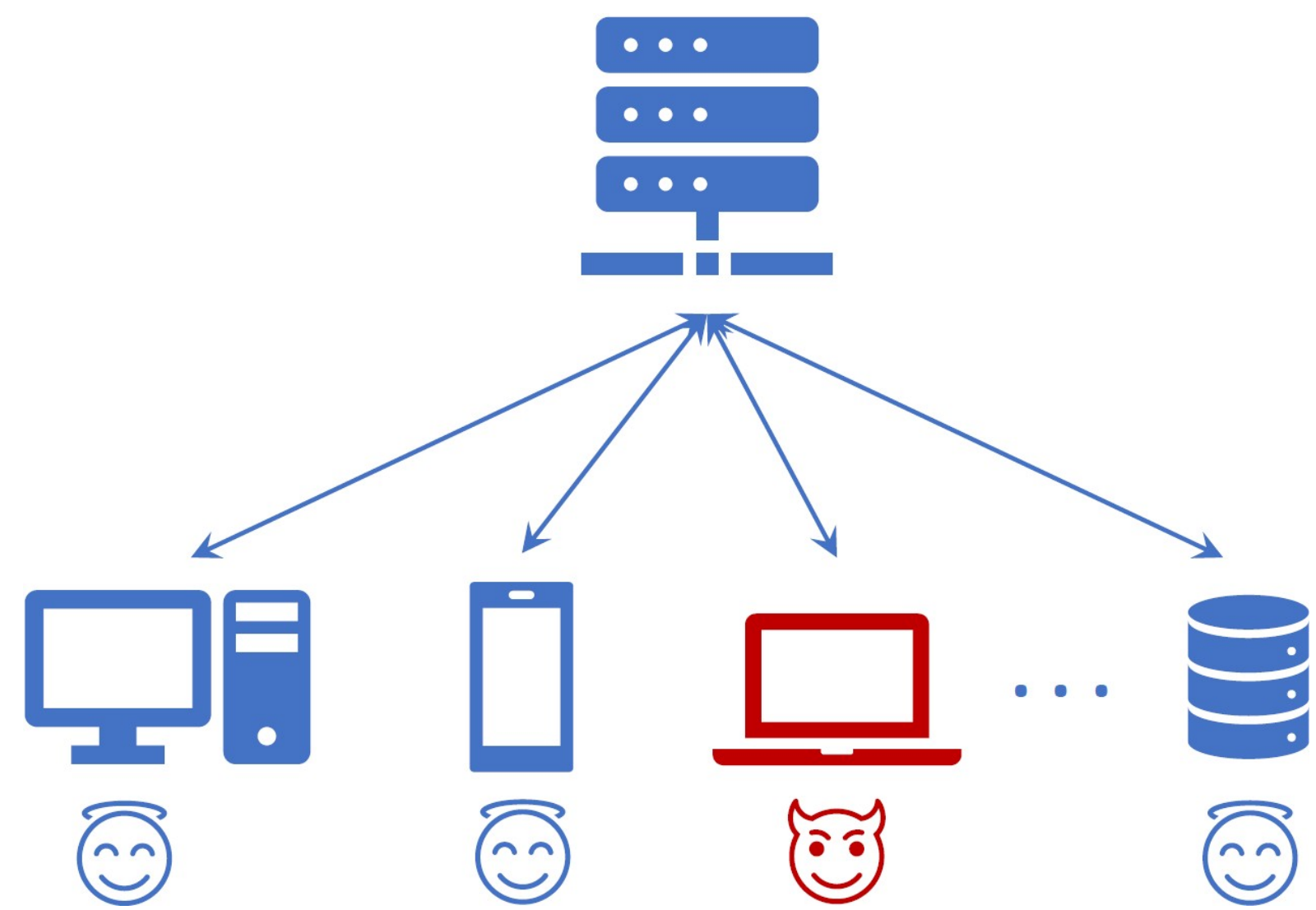


## Distributed VI problem

A lot of problems cannot be reduced to minimization, e.g., adversarial training [1], generative adversarial networks (GANs) [2], hierarchical reinforcement learning [3], adversarial examples games [4], problems arising in game theory, control theory, and differential equations [5]. Such problems lead to min-max or, more generally, variational inequality (VI) problems [6] that have significant differences from minimization ones (but include minimization) and require special consideration [7, 8].

$$\text{Find } \mathbf{x}^* \in \mathbb{R}^d \text{ s.t. } F(\mathbf{x}^*) = 0, \text{ where } F(\mathbf{x}) := \frac{1}{G} \sum_{i \in \mathcal{G}} F_i(\mathbf{x}),$$

- $\mathcal{G}$  is the set of **good clients**
- $\mathcal{B}$  is the set of **Byzantine workers** – the workers that can arbitrarily deviate from the prescribed protocol (maliciously or not) and are assumed to be omniscient
- $\mathcal{G} \cup \mathcal{B} = [n]$  is the set of clients participating in training



## Robust Aggregation

### Popular aggregation rules:

- **Krum**( $x_1, \dots, x_n$ )  $\stackrel{\text{def}}{=} \arg \min_{x_i \in \{x_1, \dots, x_n\}} \sum_{j \in S_i} \|x_j - x_i\|^2$ , where  $S_i \subseteq \{x_1, \dots, x_n\}$  are  $n - |\mathcal{B}| - 2$  closest vectors to  $x_i$
- **Robust Fed. Averaging**:  $\text{RFA}(x_1, \dots, x_n) \stackrel{\text{def}}{=} \arg \min_{x \in \mathbb{R}^d} \sum_{i=1}^n \|x - x_i\|$
- **Coordinate-wise Median**:  $[\text{CM}(x_1, \dots, x_n)]_l \stackrel{\text{def}}{=} \arg \min_{u \in \mathbb{R}} \sum_{i=1}^n |u - [x_i]_l|$

These defenses are vulnerable to Byzantine attacks [9, 10] and do not satisfy the following definition.

**Definition 1:**  $(\delta, c)$ -Robust Aggregator (modification of the definition from [11])

If a subset  $\mathcal{G} \subseteq [n]$  of  $\{x_1, x_2, \dots, x_n\}$  is s.t.  $|\mathcal{G}| = G \geq (1 - \delta)n$  for  $\delta < 0.5$  and there exists  $\sigma \geq 0$  such that  $\frac{1}{G(G-1)} \sum_{i, l \in \mathcal{G}} \mathbb{E} \|x_i - x_l\|^2 \leq \sigma^2$  where the expectation is taken w.r.t. the randomness of  $\{x_i\}_{i \in \mathcal{G}}$ , then  $\hat{x} = \text{RAgg}(x_1, \dots, x_n)$  is called  **$(\delta, c)$ -Robust Aggregator** ( $(\delta, c)$ -RAgg) if the following holds:

$$\mathbb{E} [\|\hat{x} - \bar{x}\|^2] \leq c\delta\sigma^2, \quad (1)$$

where  $\bar{x} = \frac{1}{|\mathcal{G}|} \sum_{i \in \mathcal{G}} x_i$ . If additionally  $\hat{x}$  is computed without the knowledge of  $\sigma^2$ , we say that  $\hat{x}$  is  **$(\delta, c)$ -Agnostic Robust Aggregator** ( $(\delta, c)$ -ARAgn) and write  $\hat{x} = \text{ARAgn}(x_1, \dots, x_n)$ .

One can robustify **Krum**, **RFA**, and **CM** using bucketing [11].

**Algorithm 1 Bucketing:** Robust Aggregation using bucketing [11]

1. **Input:**  $\{x_1, \dots, x_n\}$ ,  $s \in \mathbb{N}$  – bucket size, **Aggr** – aggregation rule
2. Sample random permutation  $\pi = (\pi(1), \dots, \pi(n))$  of  $[n]$
3. Compute  $y_i = \frac{1}{s} \sum_{k=s(i-1)+1}^{\min\{si, n\}} x_{\pi(k)}$  for  $i = 1, \dots, \lceil n/s \rceil$
4. **Return:**  $\hat{x} = \text{Aggr}(y_1, \dots, y_{\lceil n/s \rceil})$

## Existing Methods

**Parallel SGDA / Parallel SEG:**

- ✗ Permutation invariance
- ✗ Divergence with Byzantines

**RDEG** [12]:

- ✗ Permutation invariance
- ✗ Convergence with *large batches* in **homogeneous** case

**Why permutation non-invariance?** As [13] prove, any permutation-invariant algorithm fails to converge to any predefined accuracy even if workers have homogeneous data!

### Main Contribution

- **Methods with provably robust aggregation.** We propose new methods **SGDA-RA** and **SEG-RA** – variants of popular **SGDA** and **SEG**. We prove that **SGDA-RA** and **SEG-RA** work with any  $(\delta, c)$ -robust aggregation rule and converge to the desired accuracy *if the batchsize is large enough*.

- **Client momentum.** We add client momentum to **SGDA-RA** and propose Momentum **SGDA-RA** (**M-SGDA-RA**). That breaks the permutation invariance. In the case of star-cocoercive quasi-strongly monotone VIs, we prove the convergence to the neighborhood of the solution; the size of the neighborhood can be reduced via increasing batchsize only – similarly to the results for **RDEG**, **SGDA-RA**, and **SEG-RA**.

- **Methods with random checks of computations.** For homogeneous data case ( $\zeta = 0$ ), we propose a version of **SGDA** and **SEG** with random checks of computations (**SGDA-CC**, **SEG-CC** and their restarted versions – **R-SGDA-CC** and **R-SEG-CC**). We prove that the proposed methods converge *to any accuracy of the solution without any assumptions on the batchsize*. Moreover, when the target accuracy of the solution is small enough, the obtained convergence rates for **R-SGDA-CC** and **R-SEG-CC** are not worse than the ones for distributed **SGDA** and **SEG** derived in the case of no Byzantine workers; see the comparison of the convergence rates in Table 1.

## Methods with Robust Aggregation

$$\mathbb{E} \mathbf{g}_i(\mathbf{x}, \xi_i) = F_i(\mathbf{x}) \quad \mathbb{E} \xi_i \| \mathbf{g}_i(\mathbf{x}, \xi_i) - F_i(\mathbf{x}) \|^2 \leq \sigma^2. \quad (2)$$

$$\text{SGDA-RA:} \quad \mathbf{x}^{t+1} = \mathbf{x}^t - \gamma \text{RAGG}(\mathbf{g}_1^t, \dots, \mathbf{g}_n^t),$$

where  $\mathbf{g}_i^t = \mathbf{g}_i(\mathbf{x}^t, \xi_i^t) \forall i \in \mathcal{G}$ ,  $\mathbf{g}_i^t = * \forall i \in \mathcal{B}$ , and  $\{\mathbf{g}_i^t\}_{i \in \mathcal{G}}$  are sampled independently.

- ✗ Permutation non-invariance
- ✓ Convergence with large batches in **heterogeneous** case
- ✗ Convergence with large batches in **homogeneous** case

$$\text{M-SGDA-RA:} \quad \mathbf{x}^{t+1} = \mathbf{x}^t - \gamma \text{RAGG}(\mathbf{m}_1^t, \dots, \mathbf{m}_n^t),$$

$$\text{with } \mathbf{m}_i^t = (1 - \alpha)\mathbf{m}_i^{t-1} + \alpha\mathbf{g}_i^t,$$

where  $\mathbf{g}_i^t = \mathbf{g}_i(\mathbf{x}^t, \xi_i^t)$ ,  $\forall i \in \mathcal{G}$  and  $\mathbf{g}_i^t = * \forall i \in \mathcal{B}$  and  $\{\mathbf{g}_i^t\}_{i \in \mathcal{G}}$  are sampled independently.

- ✓ Permutation non-invariance
- ✓ Convergence with large batches in **heterogeneous** case
- ✗ Convergence with large batches in **homogeneous** case

## Methods with Checks of Computations

### Key idea of the checks

At each iteration of **SGDA-CC**, the server selects  $m$  workers (uniformly at random) and requests them to check the computations of other  $m$  workers from the previous iteration.

Let  $V_t$  be the set of workers that verify/check computations,  $A_t$  are active workers at iteration  $t$ , and  $V_t \cap A_t = \emptyset$ . Then, the update of **SGDA-CC** can be written as

$$\text{SGDA-CC:} \quad \mathbf{x}^{t+1} = \mathbf{x}^t - \gamma \bar{\mathbf{g}}^t,$$

$$\text{if } \bar{\mathbf{g}}^t = \frac{1}{|A_t|} \sum_{i \in A_t} \mathbf{g}_i(\mathbf{x}^t, \xi_i^t) \text{ is accepted,}$$

where  $\{\mathbf{g}_i(\mathbf{x}^t, \xi_i^t)\}_{i \in \mathcal{G}}$  are sampled independently.

The acceptance (of the update) event occurs when the condition  $\|\bar{\mathbf{g}}^t - \mathbf{g}_i(\mathbf{x}^t, \xi_i^t)\| \leq C\sigma$  holds for the majority of workers. If  $\bar{\mathbf{g}}^t$  is rejected, then all workers re-sample  $\mathbf{g}_i(\mathbf{x}^t, \xi_i^t)$  until acceptance is achieved. The rejection probability is bounded, as per [14].

- ✗ Permutation invariance
- ✗ Non applicable for **heterogeneous** case
- ✓ Convergence with *any batches* in **homogeneous** case

$$\text{R-SGDA-CC:} \quad \text{restarted version of SGDA-CC}$$

- ✓ Additionally benefits of collaboration

## References

- [1] Ian J Goodfellow, Jonathon Shlens, and Christian Szegedy. “Explaining and harnessing adversarial examples”. In: *International Conference on Learning Representations* (2015).
- [2] Ian Goodfellow et al. “Generative Adversarial Nets”. In: *Advances in Neural Information Processing Systems* 27 (2014).
- [3] Greg Wayne and LF Abbott. “Hierarchical control using networks trained with higher-level forward models”. In: *Neural computation* 26.10 (2014), pp. 2163–2193.
- [4] Joey Bose et al. “Adversarial example games”. In: *Advances in neural information processing systems* 33 (2020), pp. 8921–8934.
- [5] Francisco Facchinei and Jong-Shi Pang. *Finite-dimensional variational inequalities and complementarity problems*. Springer, 2003.
- [6] Gauthier Gidel et al. “A Variational Inequality Perspective on Generative Adversarial Networks”. In: *International Conference on Learning Representations*. 2018.
- [7] Patrick T Harker and Jong-Shi Pang. “Finite-dimensional variational inequality and nonlinear complementarity problems: a survey of theory, algorithms and applications”. In: *Mathematical programming* 48.1-3 (1990), pp. 161–220.
- [8] Ernest K Ryu and Wotao Yin. *Large-Scale Convex Optimization: Algorithms & Analyses via Monotone Operators*. Cambridge University Press, 2022.
- [9] Moran Baruch, Gilad Baruch, and Yoav Goldberg. *A Little Is Enough: Circumventing Defenses For Distributed Learning*. 2019. arXiv: 1902.06156 [cs.LG].
- [10] Cong Xie, Sammi Koyejo, and Indranil Gupta. *Fall of Empires: Breaking Byzantine-tolerant SGD by Inner Product Manipulation*. 2019. arXiv: 1903.03936 [cs.LG].
- [11] Sai Praneeth Karimireddy, Lie He, and Martin Jaggi. “Byzantine-Robust Learning on Heterogeneous Datasets via Bucketing”. In: *International Conference on Learning Representations*. 2022.
- [12] Arman Adibi et al. “Distributed statistical min-max learning in the presence of Byzantine agents”. In: *2022 IEEE 61st Conference on Decision and Control (CDC)*. IEEE, 2022, pp. 4179–4184.
- [13] Sai Praneeth Karimireddy, Lie He, and Martin Jaggi. “Learning from history for byzantine robust optimization”. In: *International Conference on Machine Learning*. PMLR, 2021, pp. 5311–5319.
- [14] Eduard Gorbunov et al. “Secure distributed training at scale”. In: *International Conference on Machine Learning*. PMLR, 2022, pp. 7679–7739.

## Rates and Comparison

Table 1: By the complexity, we mean the number of stochastic oracle calls needed for a method to guarantee that  $\mathbb{E} \|\mathbf{x}^T - \mathbf{x}^*\|^2 \leq \varepsilon$  (for **RDEG**  $\mathbf{P}\{\|\mathbf{x}^T - \mathbf{x}^*\|^2 \leq \varepsilon\} \geq 1 - \delta_{\text{RDEG}}$ ,  $\delta_{\text{RDEG}} \in (0, 1]$ ). Column “BS” indicates the minimal batch-size used for achieving the corresponding complexity. Notation:  $c, \delta$  are robust aggregator parameters;  $\alpha$  = momentum parameter;  $\beta$  = ratio of inner and outer stepsize in **SEG**-like methods;  $n$  = total numbers of peers;  $m$  = number of checking peers;  $G$  = number of peers following the protocol;  $R$  = any upper bound on  $\|\mathbf{x}^0 - \mathbf{x}^*\|$ ;  $\mu$  = quasi-strong monotonicity (QSM) parameter;  $\ell$  = star-cocoercivity (SC) parameter;  $L$  = Lipschitzness (Lip) parameter;  $\sigma^2$  = bound on the variance. The definition  $\mathbf{x}^T$  can vary; see corresponding theorems for the exact formulas.

Setup	Method	Complexity	BS
(SC), (QSM)	SGDA-RA	$\frac{\ell}{\mu} + \frac{1}{c\delta n}$	$\frac{c\delta\sigma^2}{\mu^2\varepsilon}$
	M-SGDA-RA	$\frac{\ell}{\mu\alpha^2} + \frac{1}{c\delta\alpha n}$	$\frac{c\delta\sigma^2}{\alpha^2\mu^2\varepsilon}$
	SGDA-CC	$\frac{\ell}{\mu} + \frac{\sigma^2}{\mu^2 n \varepsilon} + \frac{\sigma^2 n^2}{\mu^2 m \sqrt{\varepsilon}} + \frac{\sigma^2 n^2}{\mu^2 m \sqrt{\varepsilon}}$	1
	R-SGDA-CC	$\frac{\ell}{\mu} + \frac{\sigma^2}{n\mu\varepsilon} + \frac{n^2\sigma}{m\sqrt{\mu\varepsilon}}$	1
(Lip), (QSM)	SEG-RA	$\frac{L}{\beta\mu} + \frac{1}{\beta c\delta G} + \frac{1}{\beta}$	$\frac{c\delta\sigma^2}{\beta\mu^2\varepsilon}$
	SEG-CC	$\frac{L}{\mu} + \frac{1}{\beta} + \frac{\sigma^2}{\beta^2\mu^2 n \varepsilon} + \frac{\sigma^2 n^2}{\beta\mu^2 m \varepsilon} + \frac{\sigma^2 n^2}{\beta^2\mu^2 m \sqrt{\varepsilon}}$	1
	R-SEG-CC	$\frac{L}{\mu} + \frac{\sigma^2}{n\mu\varepsilon} + \frac{n^2\sigma}{m\sqrt{\mu\varepsilon}}$	1
(Lip), (QSM)	RDEG	$\frac{L}{\mu}$	$\frac{\sigma^2\mu^2 R^2}{L^4\varepsilon^2}$

- RDEG is only for homogeneous case ( $\zeta = 0$ )

Adversarial MNIST Error, attack = IPM

