

MindFlayer SGD: Efficient Parallel SGD in the Presence of Heterogeneous and Random Worker Compute Times

uai



جامعة الملك عبد الله
للعلوم والتقنية
King Abdullah University of
Science and Technology

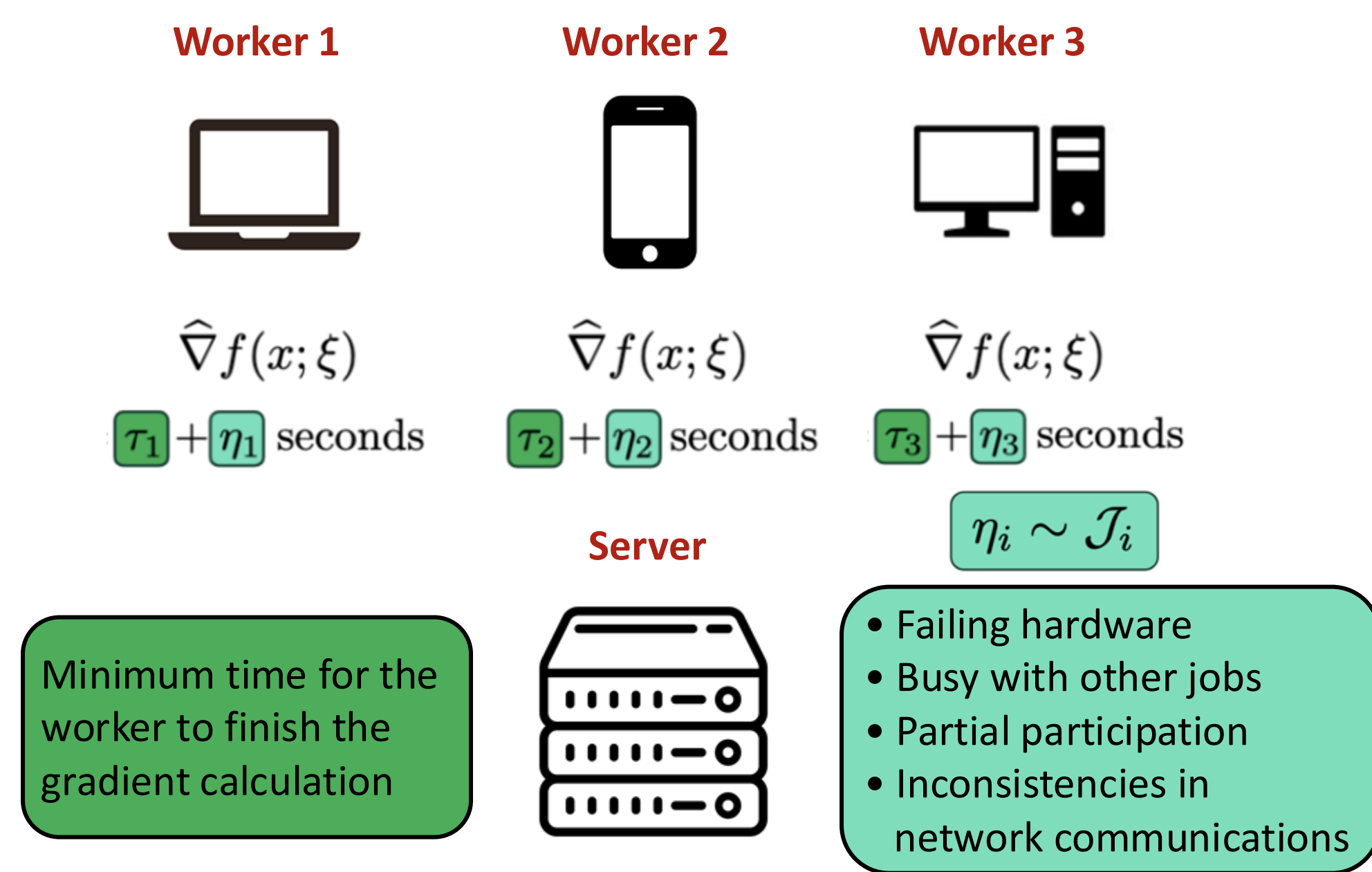
Artavazd Maranjyan, Omar Shaikh Omar, Peter Richtárik

Problem Setup: Distributed Training

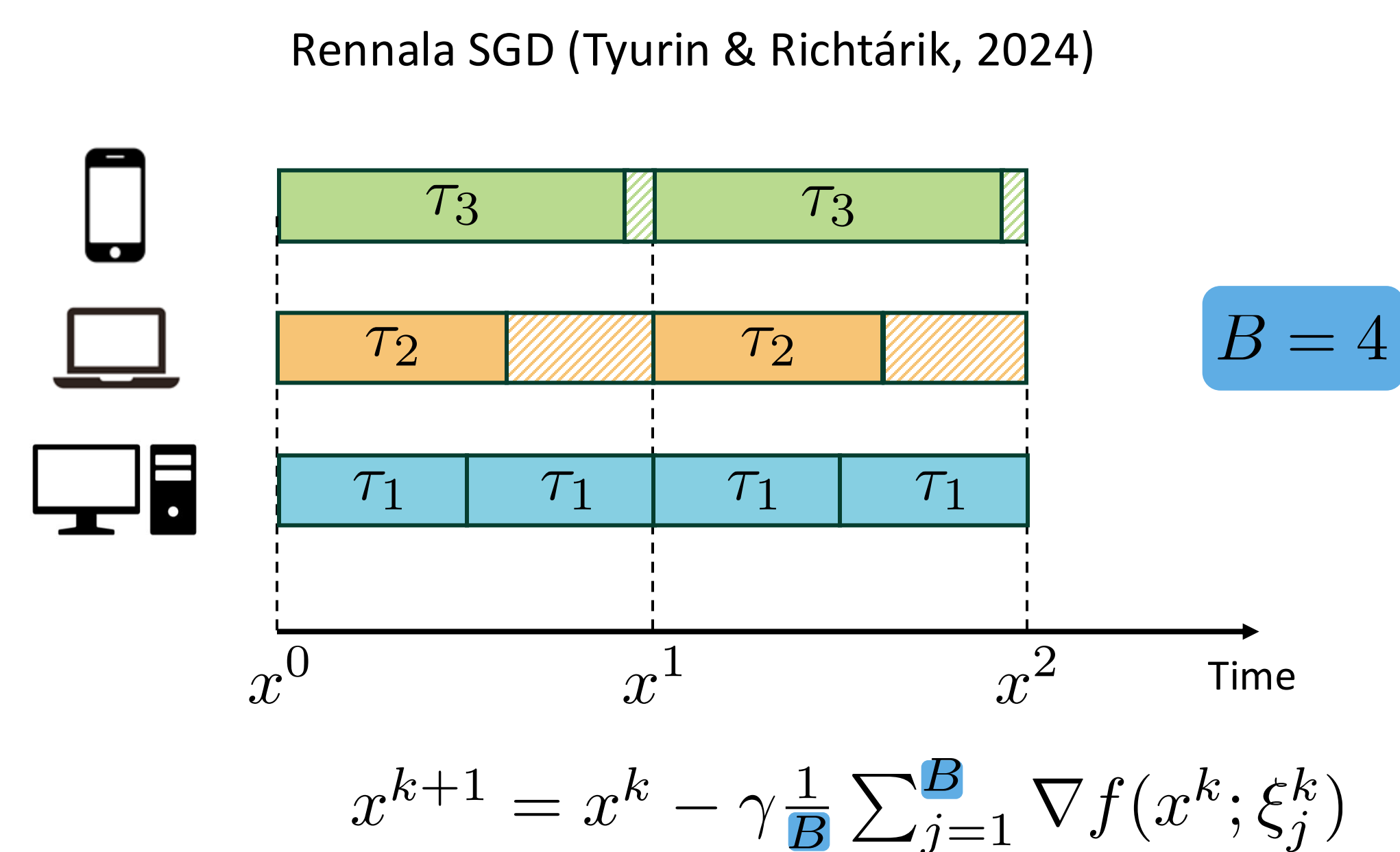
We address the nonconvex optimization problem:

$$\min_{x \in \mathbb{R}^d} \left\{ f(x) := \mathbb{E}_{\xi \sim \mathcal{D}} [f(x; \xi)] \right\},$$

We assume we have access to n parallel workers that compute stochastic gradients independently



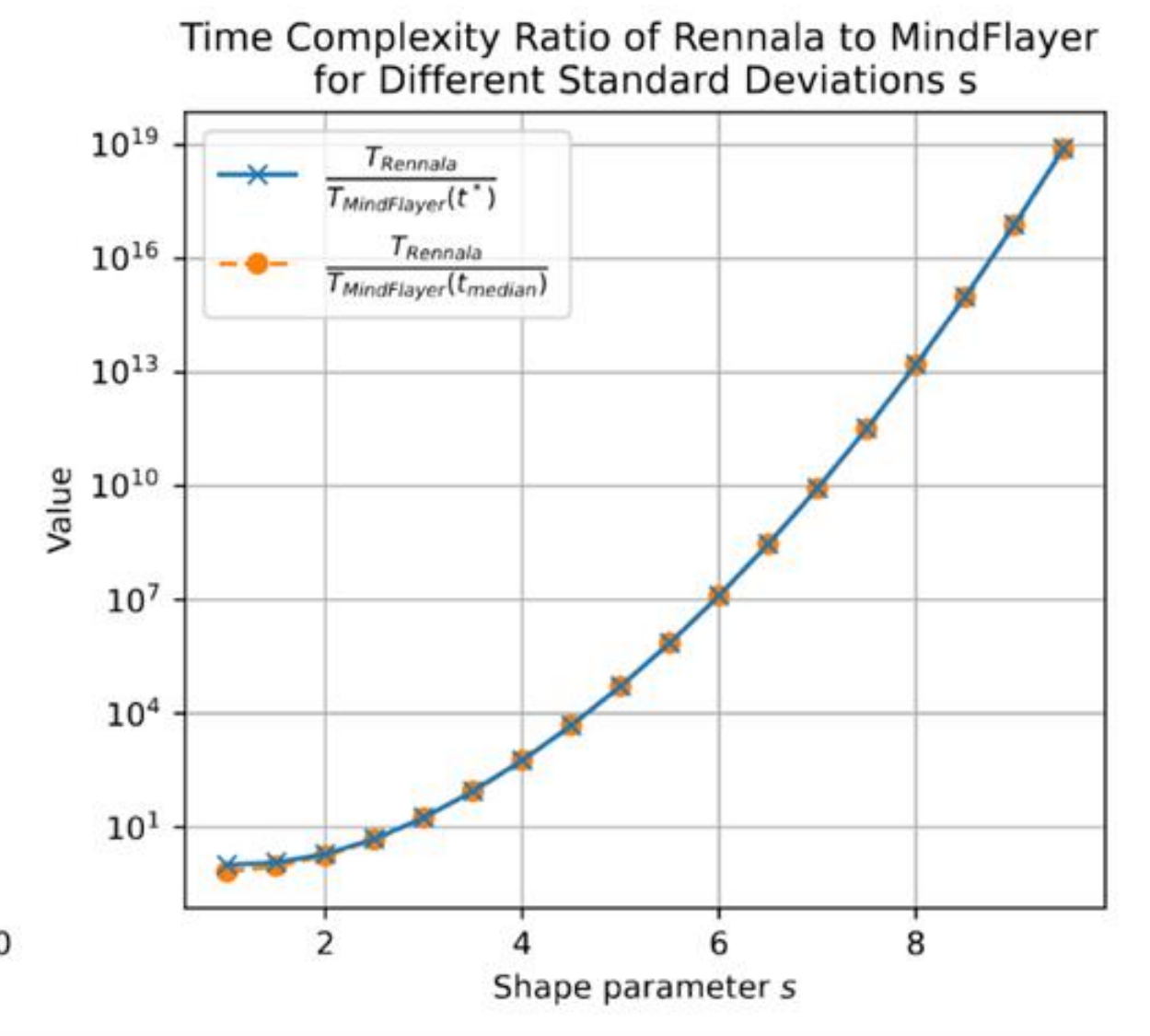
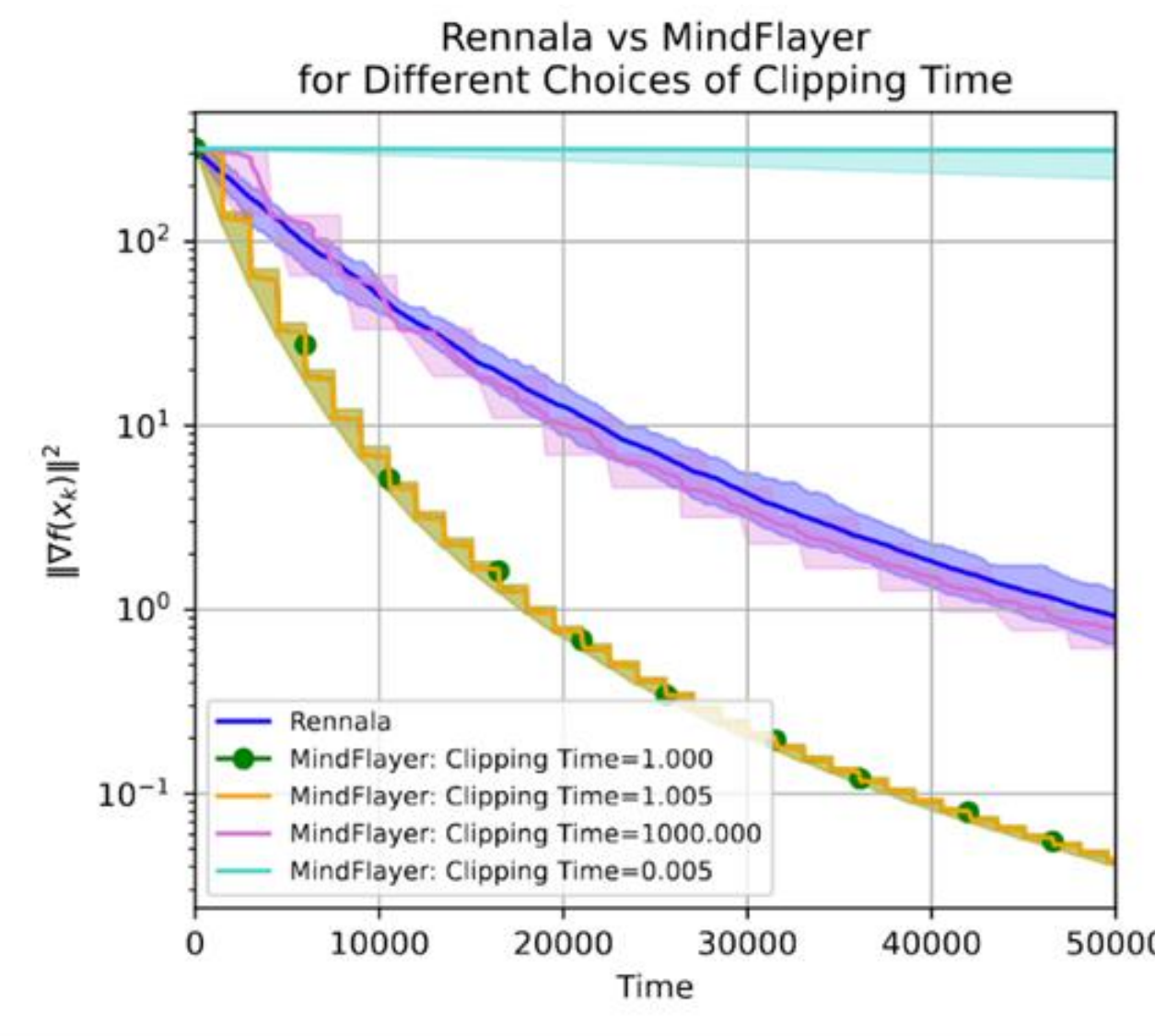
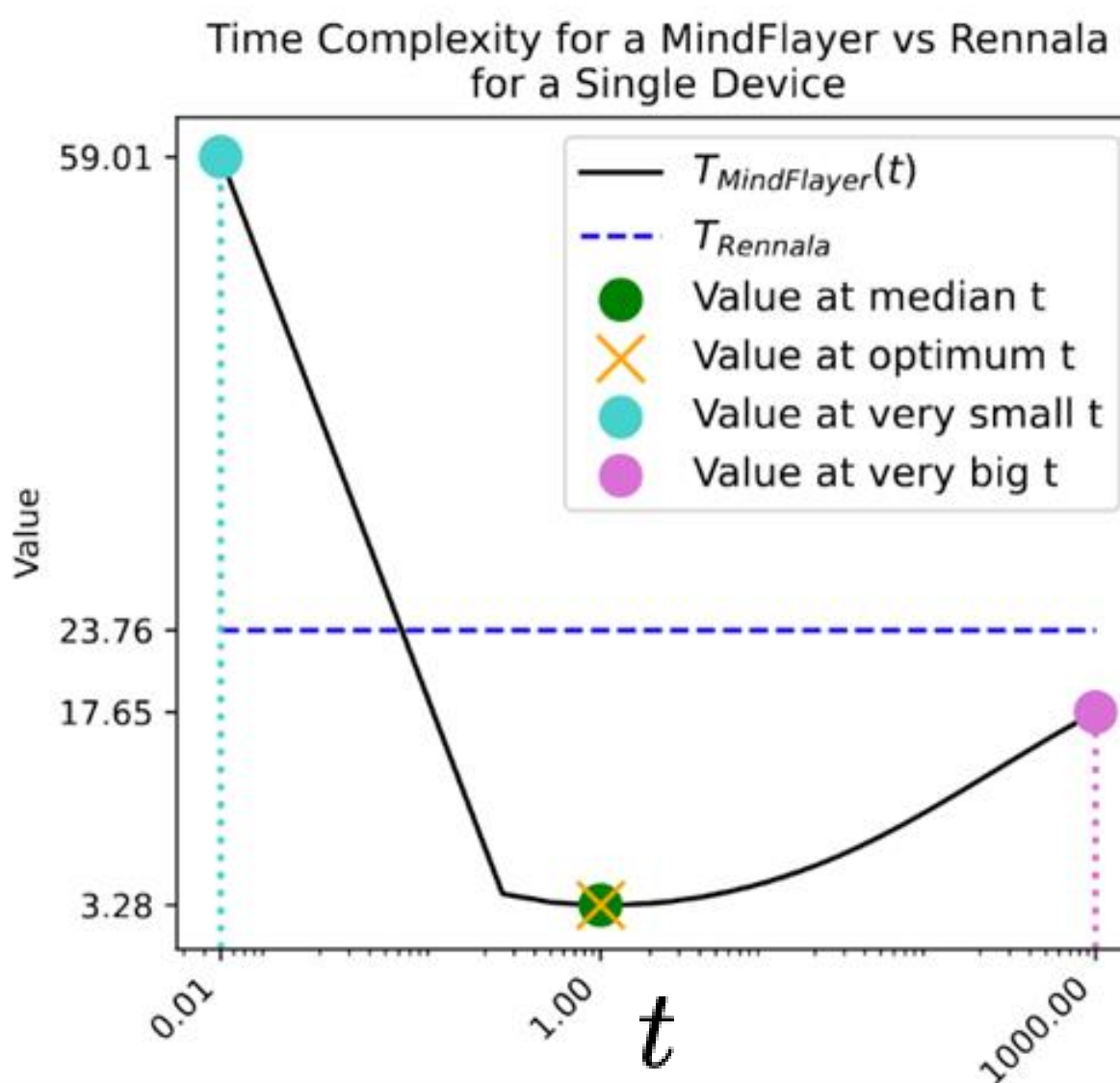
Optimal Method in Deterministic Regime



Motivation and the Single Device Case

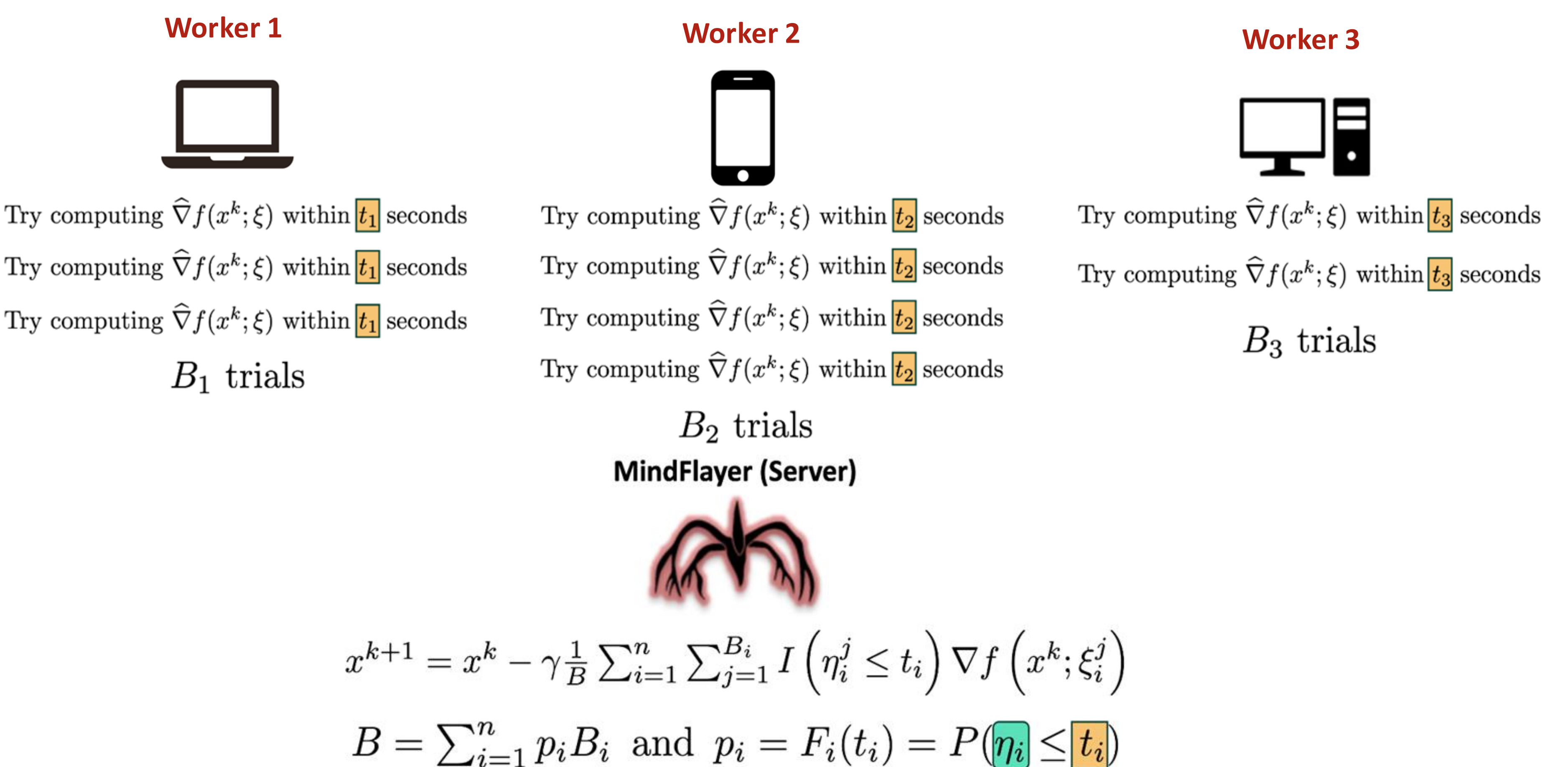
MindFlayer SGD:

- Wait t seconds for gradient response
- If the timeout expires, resend the request



$$\eta \sim \text{lognormal}(0, s)$$

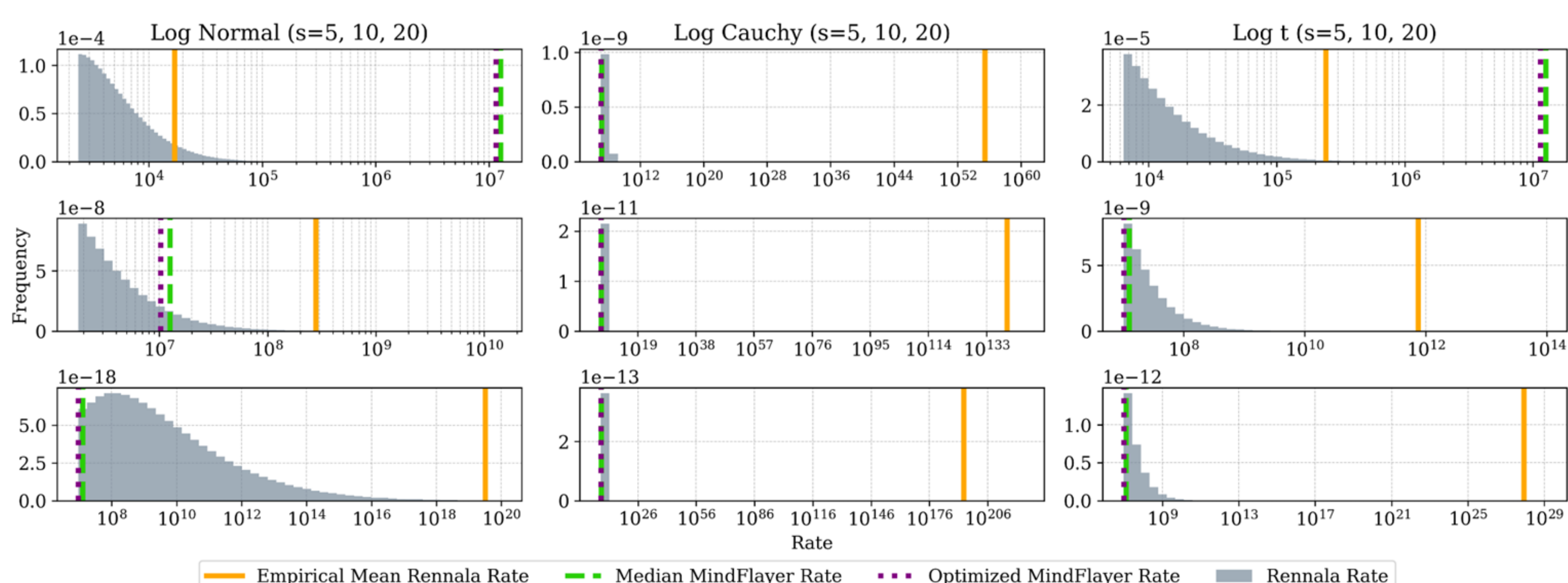
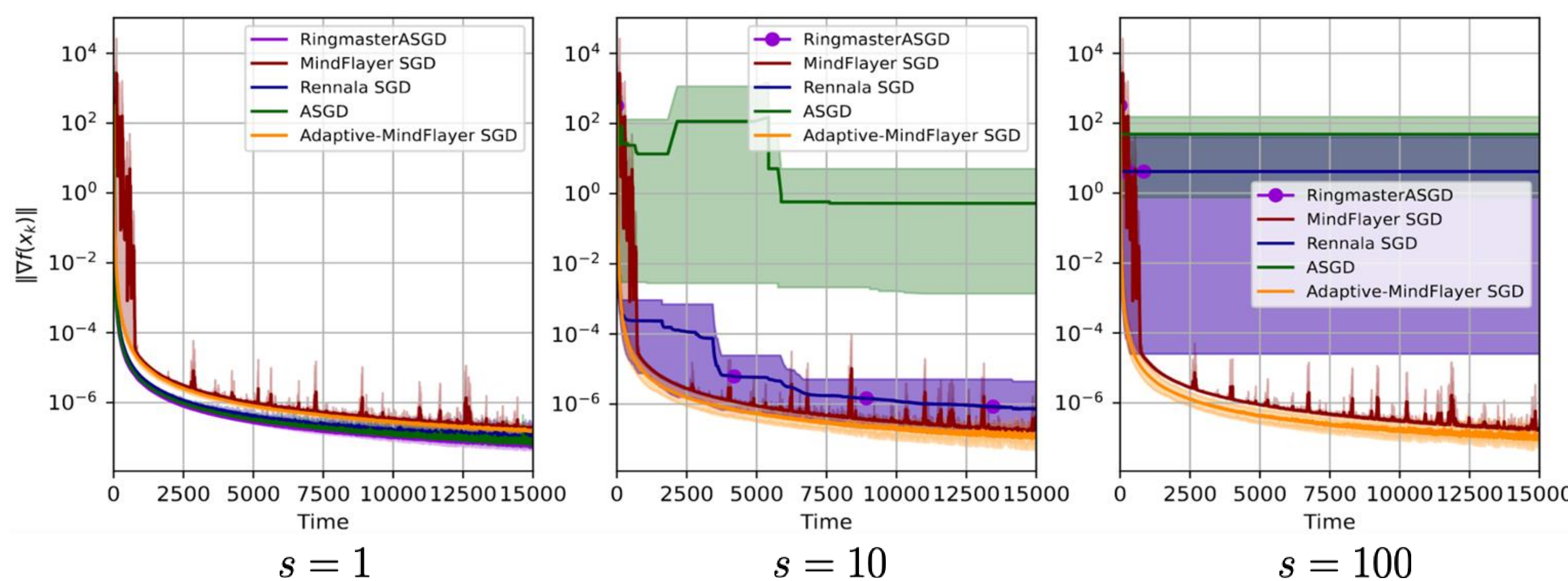
MindFlayer SGD



Experiments

$$f(x) = \frac{1}{2} x^\top A x - b^\top x \quad \forall x \in \mathbb{R}^d$$

$$\eta \sim \text{lognormal}(0, s)$$



Theoretical Results

Assumptions

Assumption 1. Function f is differentiable, and its gradient is L -Lipschitz continuous, i.e.,

$$\|\nabla f(x) - \nabla f(y)\| \leq L \|x - y\|, \text{ for all } x, y \in \mathbb{R}^d.$$

Assumption 2. The function $f(x)$ is bounded below, and we denote its infimum by $f^{\text{inf}} \in \mathbb{R}$. Let x^0 be the initial point of the optimization method, define $\Delta := f(x^0) - f^{\text{inf}}$.

Assumption 3. For all $x \in \mathbb{R}^d$, stochastic gradients $\nabla f(x; \xi)$ are unbiased and σ^2 -variance-bounded, i.e.,

$$\mathbb{E}_\xi [\nabla f(x; \xi)] = \nabla f(x),$$

$$\mathbb{E}_\xi [\|\nabla f(x; \xi) - \nabla f(x)\|^2] \leq \sigma^2,$$

where $\sigma^2 \geq 0$.

Iteration Complexity

Theorem 1. Let

$$B = \sum_{i=1}^n p_i B_i \quad \text{and} \quad \gamma = \frac{1}{2L} \min \left\{ 1, \frac{\varepsilon B}{\sigma^2} \right\}$$

Then, the method guarantees that $\frac{1}{K} \sum_{k=0}^{K-1} \mathbb{E} [\|\nabla f(x^k)\|^2] \leq \varepsilon$ if

$$K \geq \max \left\{ 1, \frac{\sigma^2}{\varepsilon B} \right\} \frac{8L (f(x^0) - f^{\text{inf}})}{\varepsilon}$$

Time Complexity

$$p_j = F_j(t_j) = P(\eta_j \leq t_j)$$

$$\Delta = f(x^0) - f^{\text{inf}}$$

$$\min_{m \in [n]} \left\{ \left(\frac{1}{m} \sum_{j=1}^m \frac{p_j}{\tau_j + t_j} \right)^{-1} \left(\frac{S}{m} + \frac{1}{m} \sum_{j=1}^m p_j \right) \frac{\Delta L}{\varepsilon} \right\}$$

$$S = \max \left\{ \frac{\sigma^2}{\varepsilon}, 1 \right\}$$