

# Controlling Delay in Asynchronous SGD: Optimality for Any Data Regime

Artavazd Maranjyan

<https://artomaranjyan.github.io>

Protocol Learning: Decentralized Collaborative Learning at Scale

ICLR 2026 Workshop

26 April 2026

**PFL**

King Abdullah University of  
Science and Technology



جامعة الملك عبد الله  
للعلوم والتقنية

# Distributed Learning (Data Parallelism)



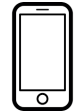
$\mathcal{D}_1$

$$f_1(x) := \mathbb{E}_{\xi_1 \sim \mathcal{D}_1} [f_1(x; \xi_1)]$$



$\mathcal{D}_2$

$$f_2(x) := \mathbb{E}_{\xi_2 \sim \mathcal{D}_2} [f_2(x; \xi_2)]$$



$\mathcal{D}_3$

$$f_3(x) := \mathbb{E}_{\xi_3 \sim \mathcal{D}_3} [f_3(x; \xi_3)]$$

Homogenous Data Setting

$$\mathcal{D}_1 = \mathcal{D}_2 = \dots = \mathcal{D}_n = \mathcal{D}$$

$$\text{minimize}_{x \in \mathbb{R}^d} \{ f(x) := \mathbb{E}_{\xi \in \mathcal{D}} [f(x; \xi)] \}$$

$$\text{minimize}_{x \in \mathbb{R}^d} \left\{ f(x) := \frac{1}{n} \sum_{i=1}^n f_i(x) \right\}$$



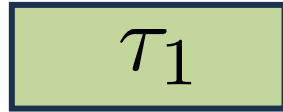
Server

# Compute Heterogeneous System



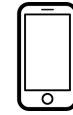
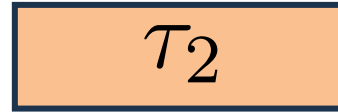
$$\nabla f_1(x; \xi_1)$$

Compute time =  $\tau_1$



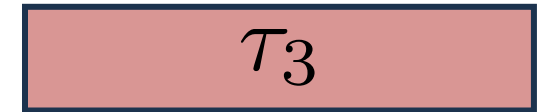
$$\nabla f_2(x; \xi_2)$$

Compute time =  $\tau_2$



$$\nabla f_3(x; \xi_3)$$

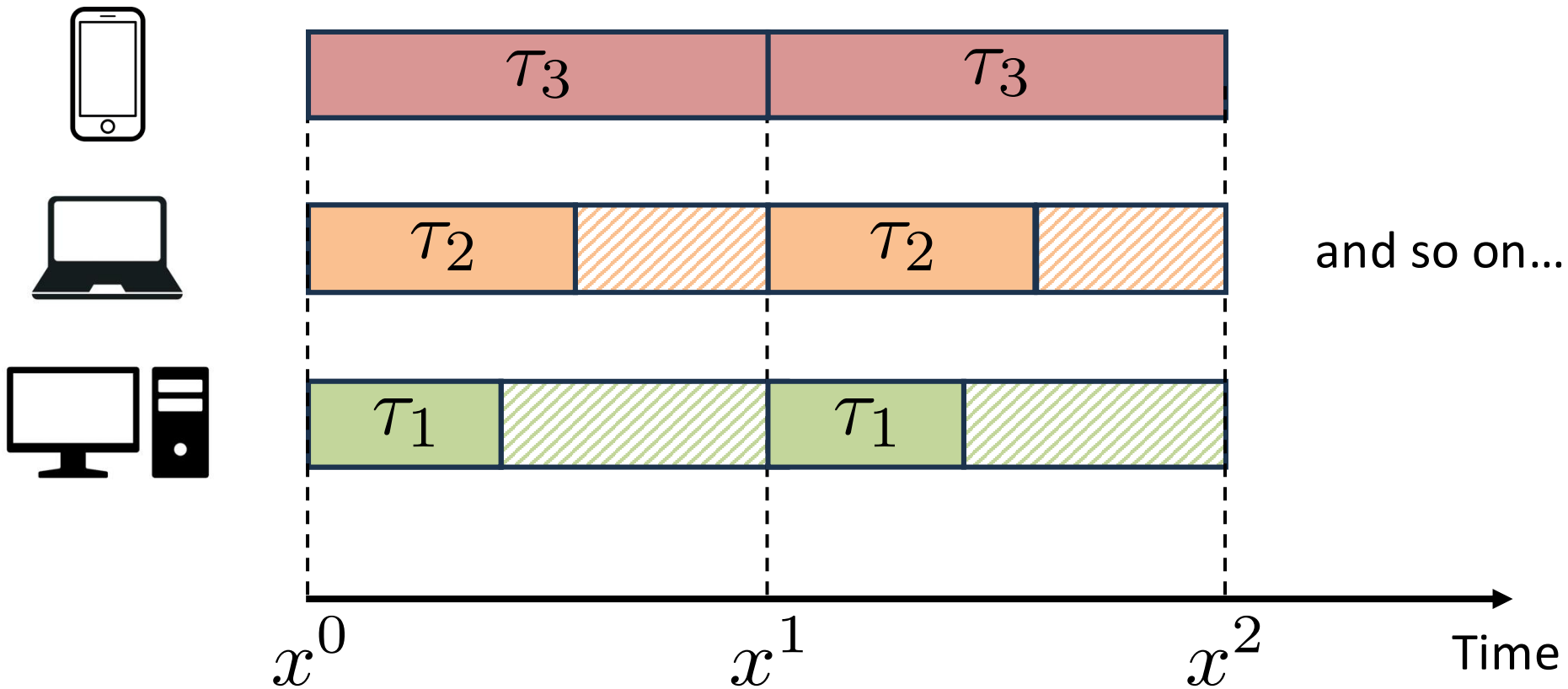
Compute time =  $\tau_3$



Server

# Naive Minibatch SGD: Each worker does one job only

$n = 3$



$$x^{k+1} = x^k - \gamma \frac{1}{n} \sum_{i=1}^n \nabla f_i(x^k; \xi_i^k)$$

# The First Asynchronous SGD



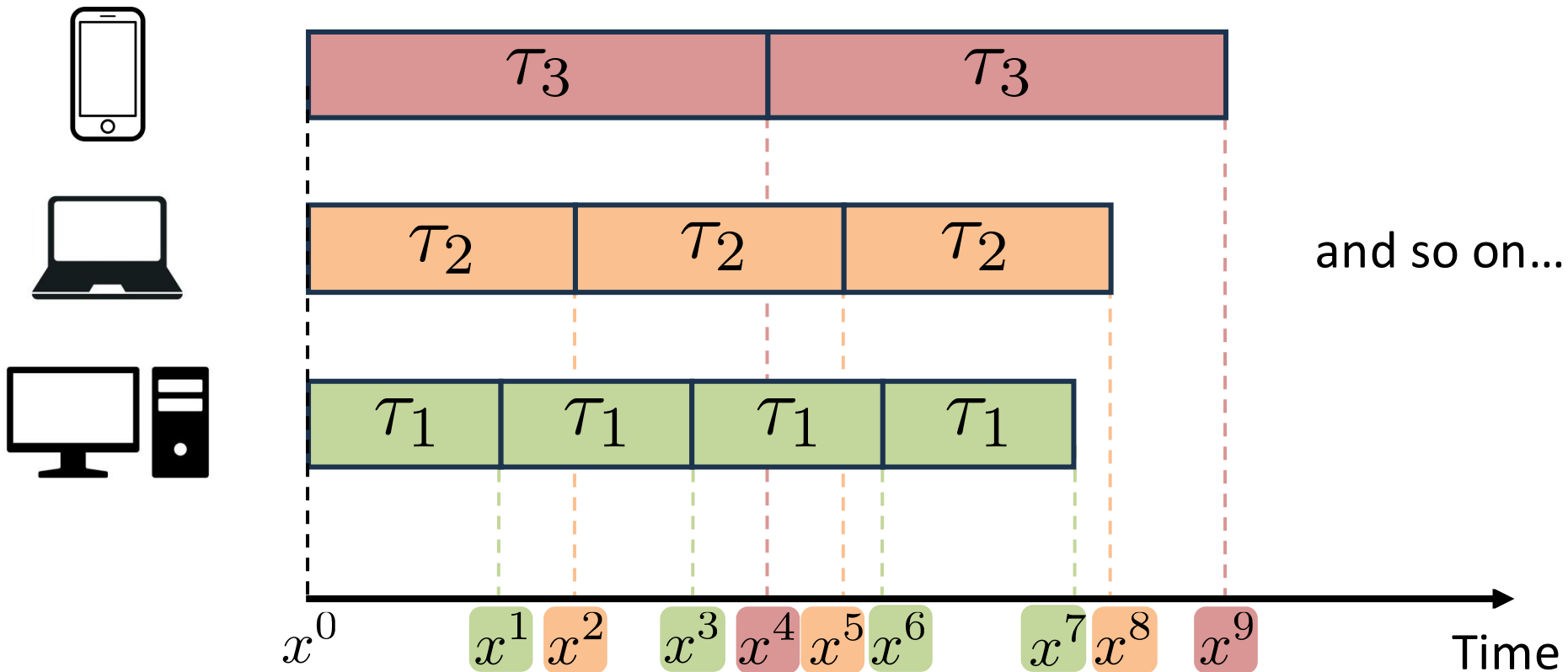
Benjamin Recht, Christopher Re, Stephen Wright, Feng Niu

**HOGWILD!: A lock-free approach to parallelizing stochastic gradient descent**

*Conference on Neural Information Processing Systems (2011)*

NeurIPS Test of Time Award, 2020

# Asynchronous SGD: Remove the synchronization



# Updates of Asynchronous SGD has delayed stochastic gradients

Delay of worker  $i^k$  at iteration  $k$

$$x^{k+1} = x^k - \gamma \nabla f_{i^k} \left( x^{k - \delta_{i^k}^k}; \xi^{k - \delta_{i^k}^k} \right)$$

Uses a gradient from a single worker



Server



# Homogeneous Data Setting

# Homogenous Data Setting

 $\mathcal{D}$ 

$$f(x) := \mathbb{E}_{\xi \sim \mathcal{D}} [f(x; \xi)]$$

 $\mathcal{D}$ 

$$f(x) := \mathbb{E}_{\xi \sim \mathcal{D}} [f(x; \xi)]$$

 $\mathcal{D}$ 

$$f(x) := \mathbb{E}_{\xi \sim \mathcal{D}} [f(x; \xi)]$$

$$\underset{x \in \mathbb{R}^d}{\text{minimize}} \{ f(x) := \mathbb{E}_{\xi \in \mathcal{D}} [f(x; \xi)] \}$$

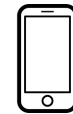


Server

# Asynchronous SGD is too wild: Ringmaster ASGD *tames* it



# Ringmaster ASGD: Have a threshold on delays



If:  $\delta^k < R$

$$x^{k+1} = x^k - \gamma \nabla f \left( x^{k-\delta^k}; \xi_i^{k-\delta^k} \right)$$

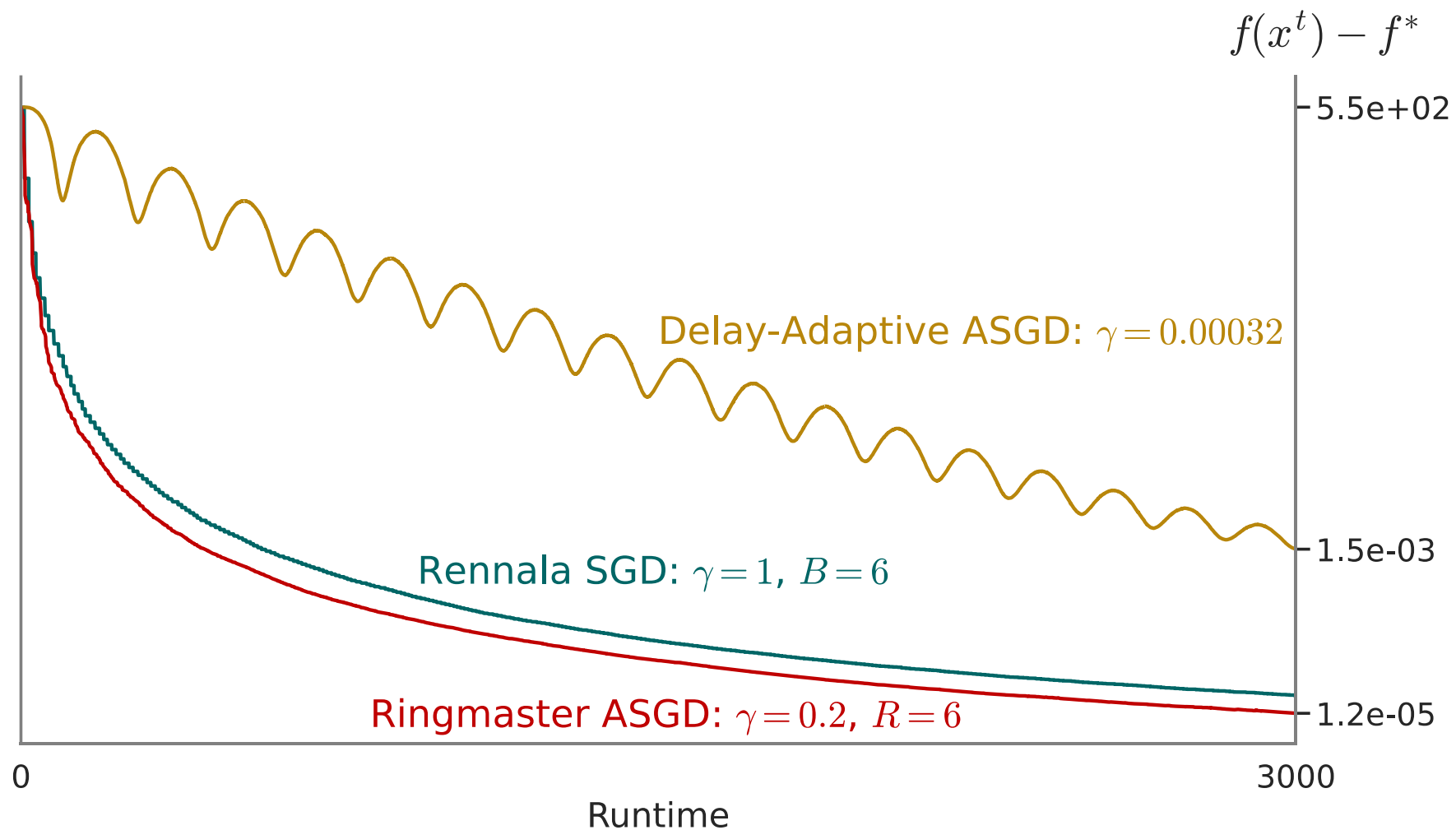
Else: Ignore the gradient and send the current point  $x^k$  to the worker



$$\nabla f \left( x^{k-\delta^k}; \xi_i^{k-\delta^k} \right)$$



# Ringmaster ASGD outperforms existing baselines





# Heterogenous Data Setting

# Heterogenous Data Setting



$\mathcal{D}_1$

$$f_1(x) := \mathbb{E}_{\xi_1 \sim \mathcal{D}_1} [f_1(x; \xi_1)]$$



$\mathcal{D}_2$

$$f_2(x) := \mathbb{E}_{\xi_2 \sim \mathcal{D}_2} [f_2(x; \xi_2)]$$



$\mathcal{D}_3$

$$f_3(x) := \mathbb{E}_{\xi_3 \sim \mathcal{D}_3} [f_3(x; \xi_3)]$$

$$\text{minimize}_{x \in \mathbb{R}^d} \left\{ f(x) := \frac{1}{n} \sum_{i=1}^n f_i(x) \right\}$$



Server

# Problems with the Naive Approach

Delay of worker  $i^k$  at iteration  $k$

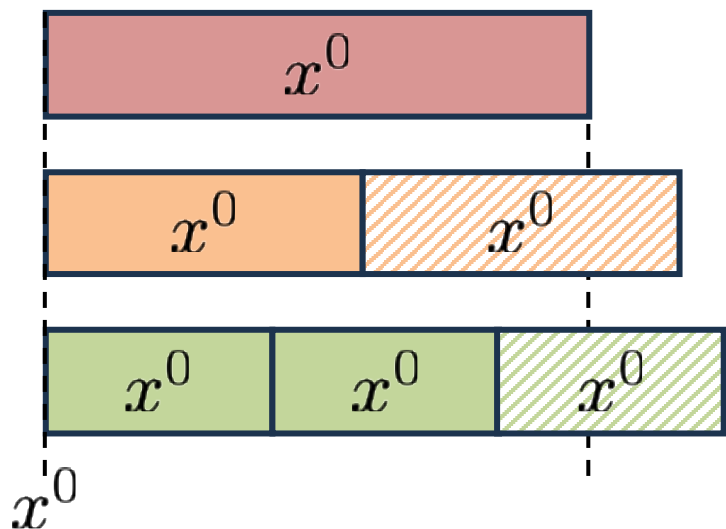
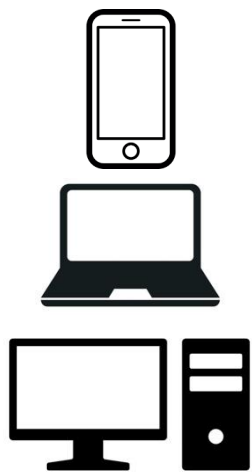
$$x^{k+1} = x^k - \gamma \nabla f_{i^k} \left( x^{k - \delta_{i^k}^k}; \xi^{k - \delta_{i^k}^k} \right)$$

Uses a gradient from a single worker



Server

# Ringleader ASGD (Phase 1)

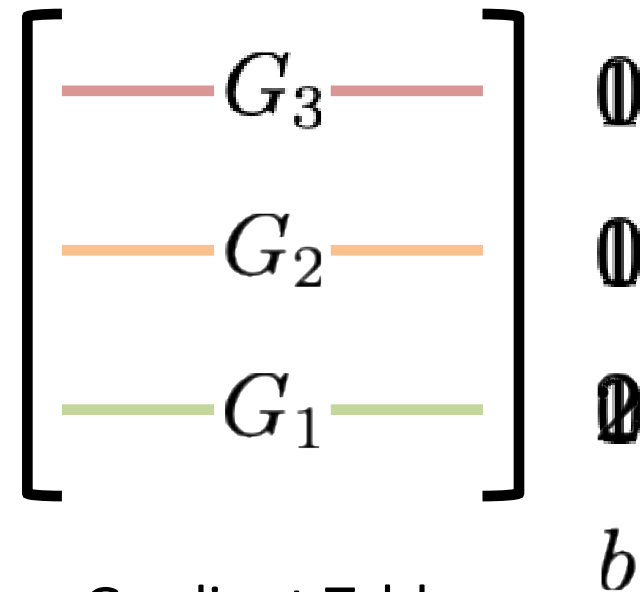


Phase 1

$$G_3 = \nabla f_3(x^0; \xi_3^{0,1})$$

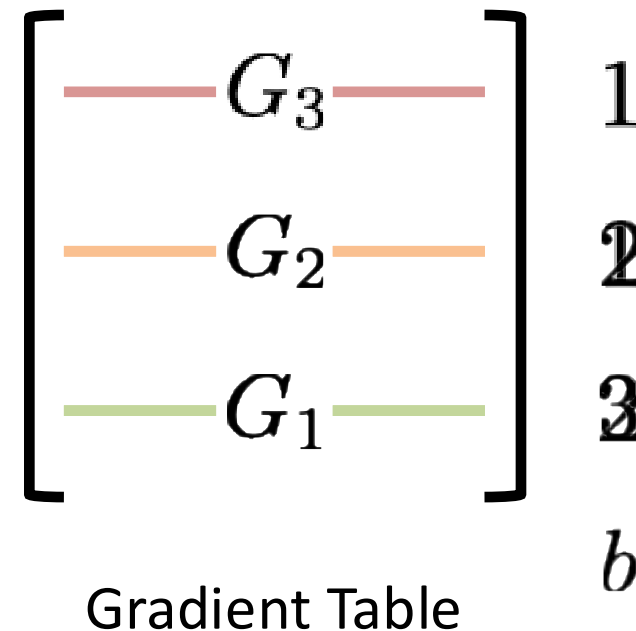
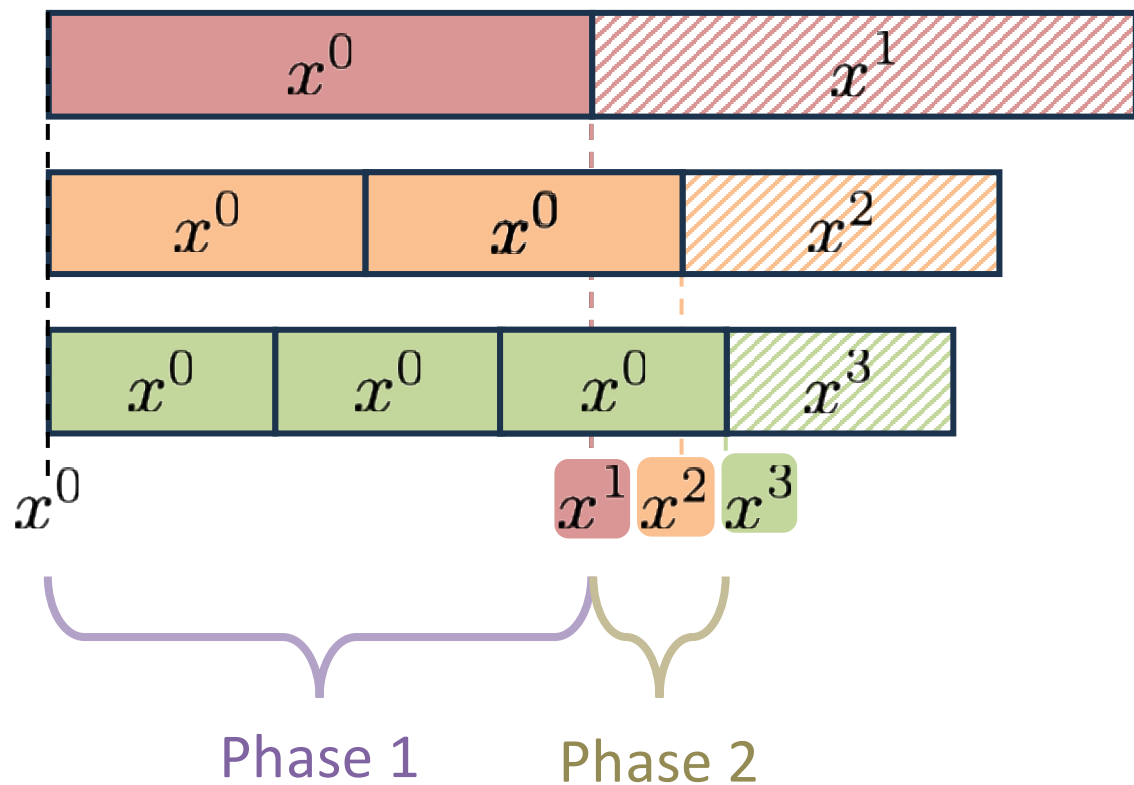
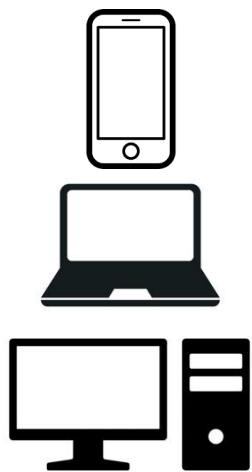
$$G_2 = \nabla f_2(x^0; \xi_2^{0,1})$$

$$G_1 = \sum_{j=1}^2 \nabla f_1(x^0; \xi_1^{0,j})$$



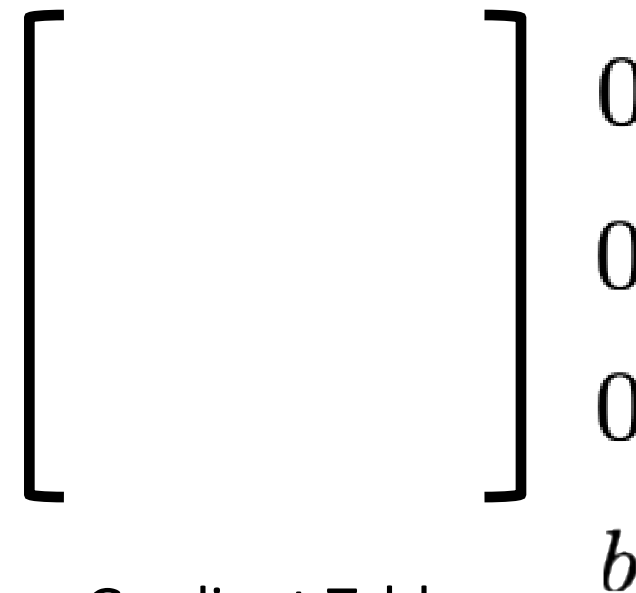
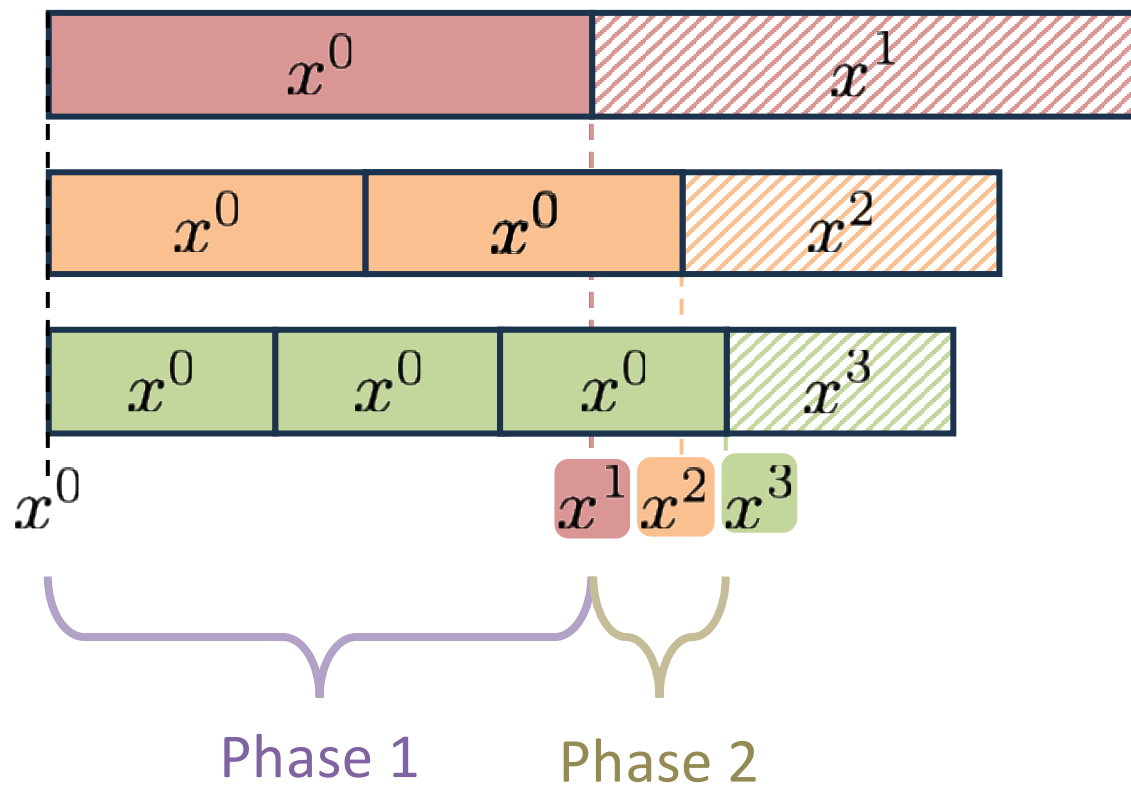
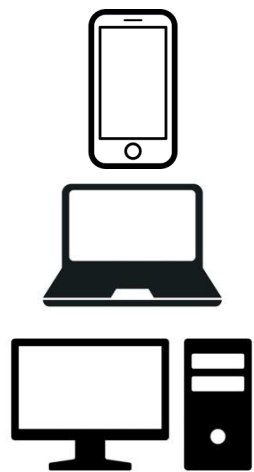
Gradient Table

# Ringleader ASGD (Phase 2)



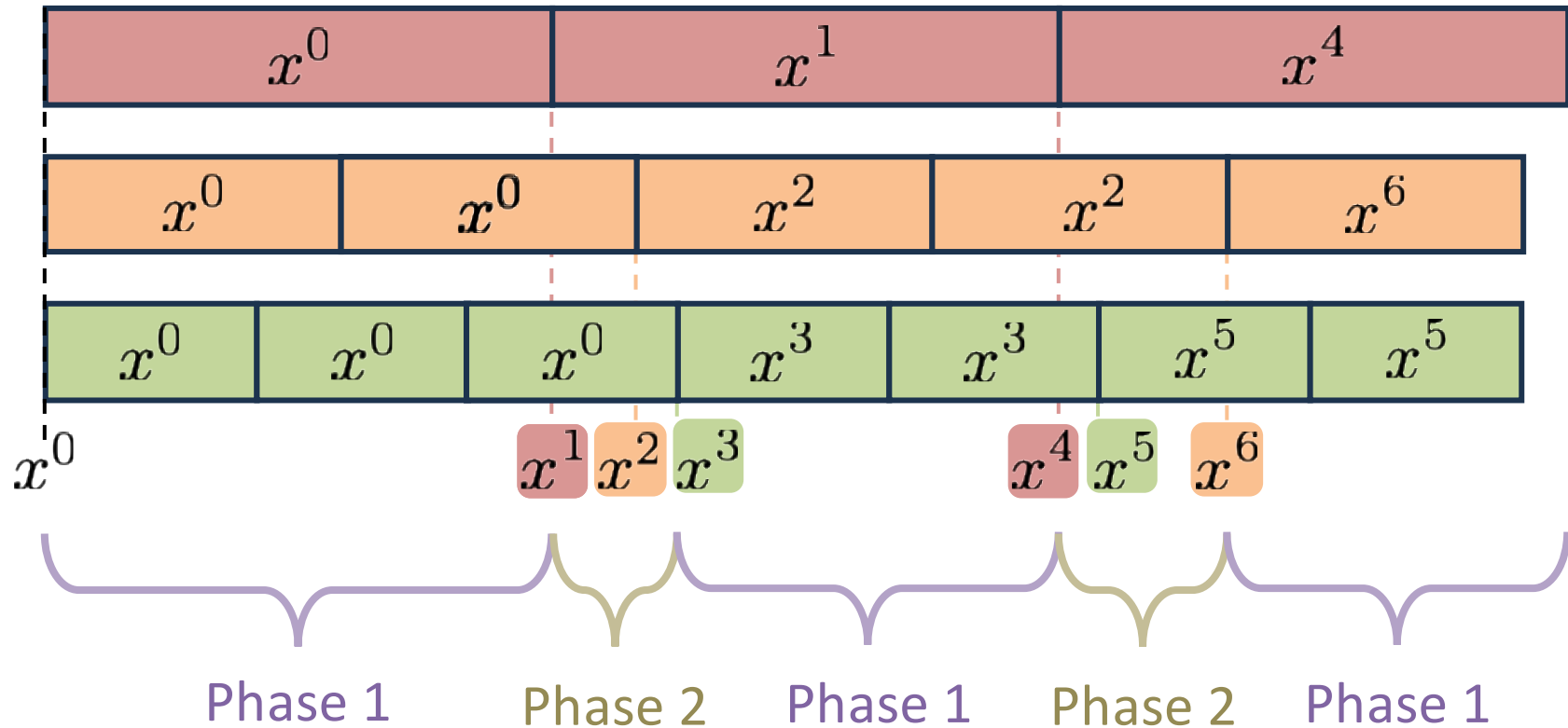
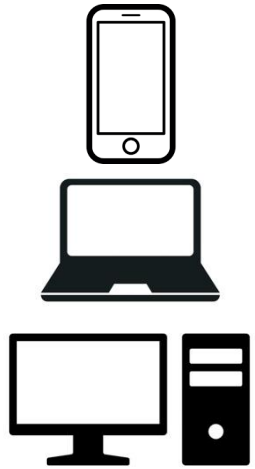
$$x^3 = x^2 - \gamma \frac{1}{n} \sum_{i=1}^n \frac{G_i}{b_i}$$

# Ringleader ASGD (Transition Phase)



Gradient Table

# Ringleader ASGD



$$x^{k+1} = x^k - \gamma \frac{1}{n} \sum_{i=1}^n \frac{G_i}{b_i}$$

$$\delta_i^k \leq 2n - 2$$

# Ringleader ASGD outperforms existing baselines

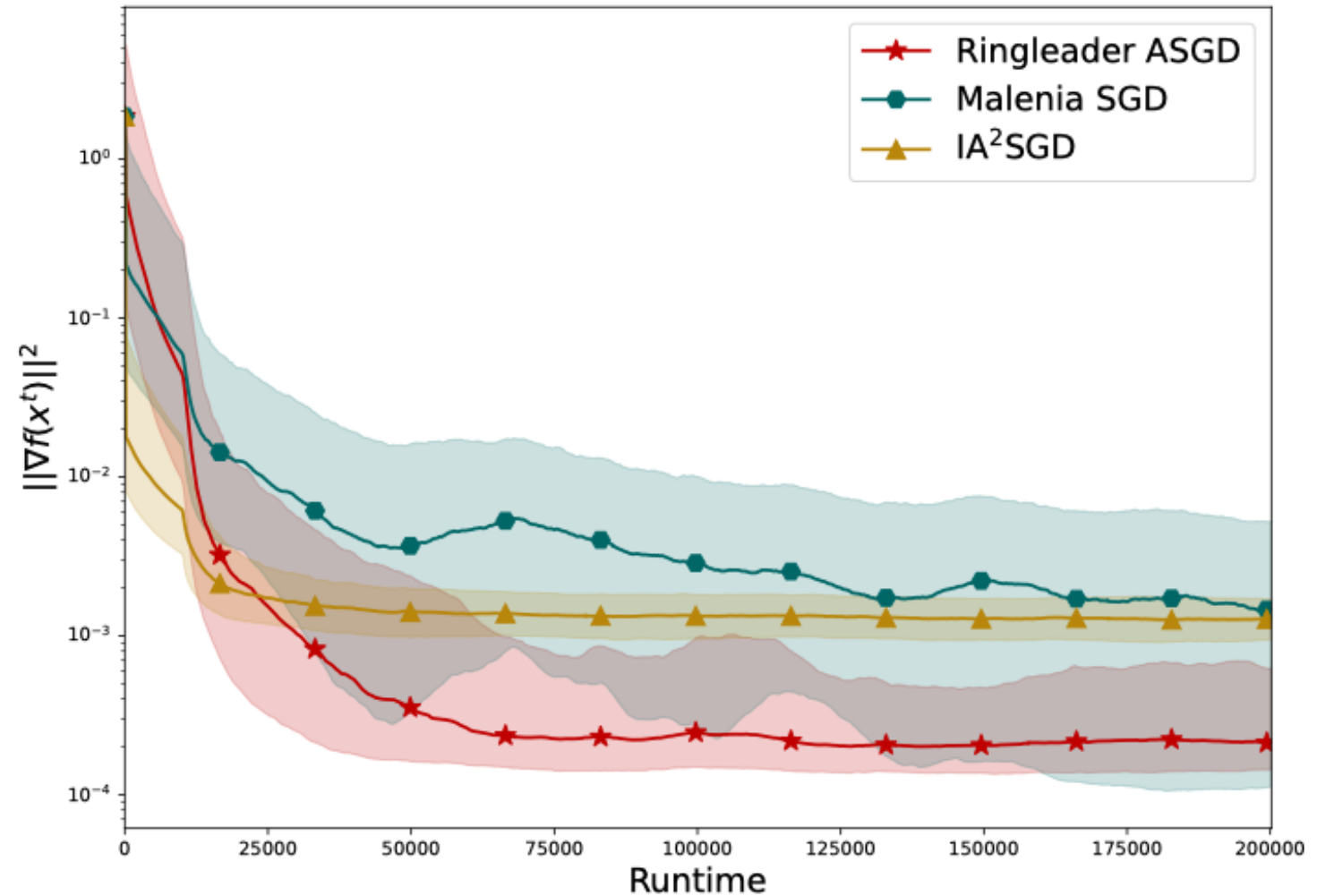
Two-layer MLP

MNIST

$n = 100$

$\tau_i = i + |\eta_i|$

$\eta_i \sim \mathcal{N}(0, i)$





**KHALAS**