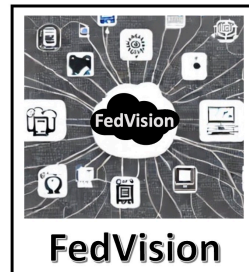# The First Optimal Parallel SGD
## (in the Presence of Data, Compute and Communication Heterogeneity)
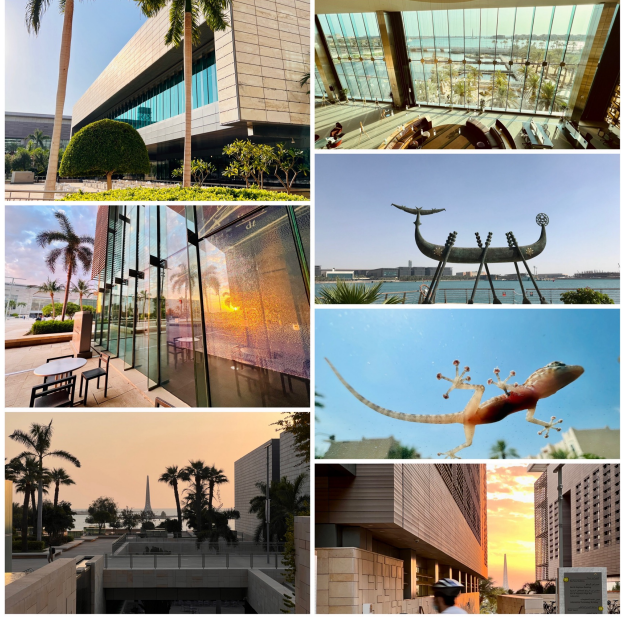
**Peter Richtárik**

King Abdullah University of Science and Technology
Kingdom of Saudi Arabia

**3rd Workshop on
Federated Learning for
Computer Vision**

in Conjunction with CVPR 2024
(6/17 All Day)

# Optimization & Machine Learning Lab @ KAUST

# Part 1
# Federated Learning

Jakub Konečný

H Brendan McMahan

Google

THE UNIVERSITY of EDINBURGH

**Federated Learning was developed in 2015/2016 in a collaboration between the University of Edinburgh & Google**

H. Brendan McMahan, Eider Moore, Daniel Ramage, Seth Hampson, Blaise Agüera y Arcas
**Communication-Efficient Learning of Deep Networks from Decentralized Data**
*20th International Conference on Artificial Intelligence and Statistics (AISTATS), 2017*

Keith Bonawitz et al
**Practical Secure Aggregation for Federated Learning on User-Held Data**
*NIPS Private Multi-Party Machine Learning Workshop, 2016*



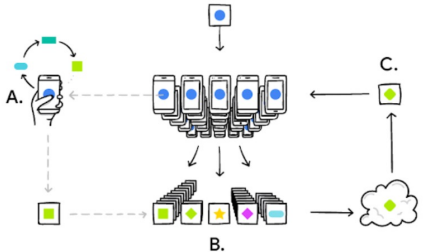Google AI Blog

The latest from Google Research

### Federated Learning: Collaborative Machine Learning without Centralized Training Data

Thursday, April 6, 2017

Posted by Brendan McMahan and Daniel Ramage, Research Scientists

Standard machine learning approaches require centralizing the training data on one machine or in a datacenter. And Google has built one of the most secure and robust cloud infrastructures for processing this data to make our services better. Now for models trained from user interaction with mobile devices, we're introducing an additional approach: *Federated Learning*.

Federated Learning enables mobile phones to collaboratively learn a shared prediction model while keeping all the training data on device, decoupling the ability to do machine learning from the need to store the data in the cloud. This goes beyond the use of local models that make predictions on mobile devices (like the Mobile Vision API and On-Device Smart Reply) by bringing model *training* to the device as well.

It works like this: your device downloads the current model, improves it by learning from data on your phone, and then summarizes the changes as a small focused update. Only this update to the model is sent to the cloud, using encrypted communication, where it is immediately averaged with other user updates to improve the shared model. All the training data remains on your device, and no individual updates are stored in the cloud.
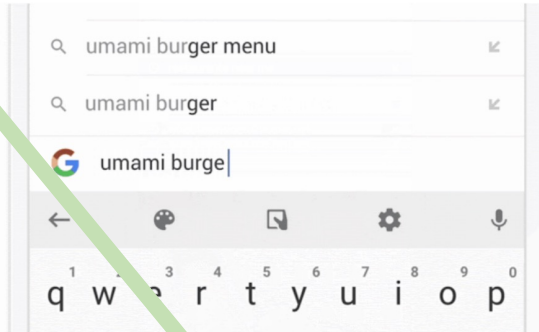
Your phone personalizes the model locally, based on your usage (A). Many users' updates are aggregated (B) to form a consensus change (C) to the shared model, after which the procedure is repeated.

Federated Learning allows for smarter models, lower latency, and less power consumption, all while ensuring privacy. And this approach has another immediate benefit: in addition to providing an update to the shared model, the improved model on your phone can also be used immediately, powering experiences personalized by the way you use your phone.

We're currently testing Federated Learning in Gboard on Android, the Google Keyboard. When Gboard shows a suggested query, your phone locally stores information about the current context and whether you clicked the suggestion. Federated Learning processes that history on-device to suggest improvements to the next iteration of Gboard's query suggestion model.
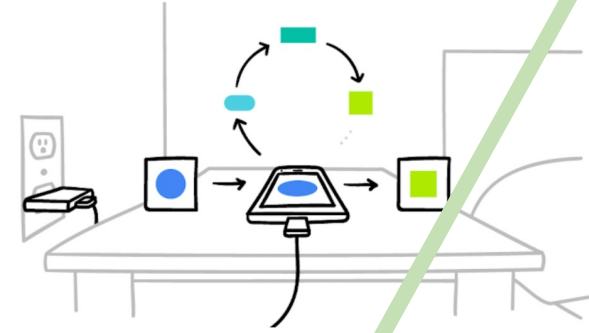
To make Federated Learning possible, we had to overcome many algorithmic and technical challenges. In a typical machine learning system, an optimization algorithm like Stochastic Gradient Descent (SGD) runs on a large dataset partitioned homogeneously across servers in the cloud. Such highly iterative algorithms require low-latency, high-throughput connections to the training data. But in the Federated Learning setting, the data is distributed across millions of devices in a highly uneven fashion. In addition, these devices have significantly higher-latency, lower-throughput connections and are only intermittently available for training.

These bandwidth and latency limitations motivate our Federated Averaging algorithm, which can train deep networks using 10-100x less communication compared to a naively federated version of SGD. The key idea is to use the powerful processors in modern mobile devices to compute higher quality updates than simple gradient steps. Since it takes fewer iterations of high-quality updates to produce a good model, training can use much less communication. As upload speeds are typically much slower than download speeds, we also developed a novel way to reduce upload communication costs up to another 100x by compressing updates using random rotations and quantization. While these approaches are focused on training deep networks, we've also designed algorithms for high-dimensional sparse convex models which excel on problems like click-through-rate prediction.

Deploying this technology to millions of heterogenous phones running Gboard requires a sophisticated technology stack. On device training uses a miniature version of TensorFlow. Careful scheduling ensures training happens only when the device is idle, plugged in, and on a free wireless connection, so there is no impact on the phone's performance.

Your phone participates in Federated Learning only when it won't negatively impact your experience.

The system then needs to communicate and aggregate the model updates in a secure, efficient, scalable, and fault-tolerant way. It's only the combination of research with this infrastructure that makes the benefits of Federated Learning possible.
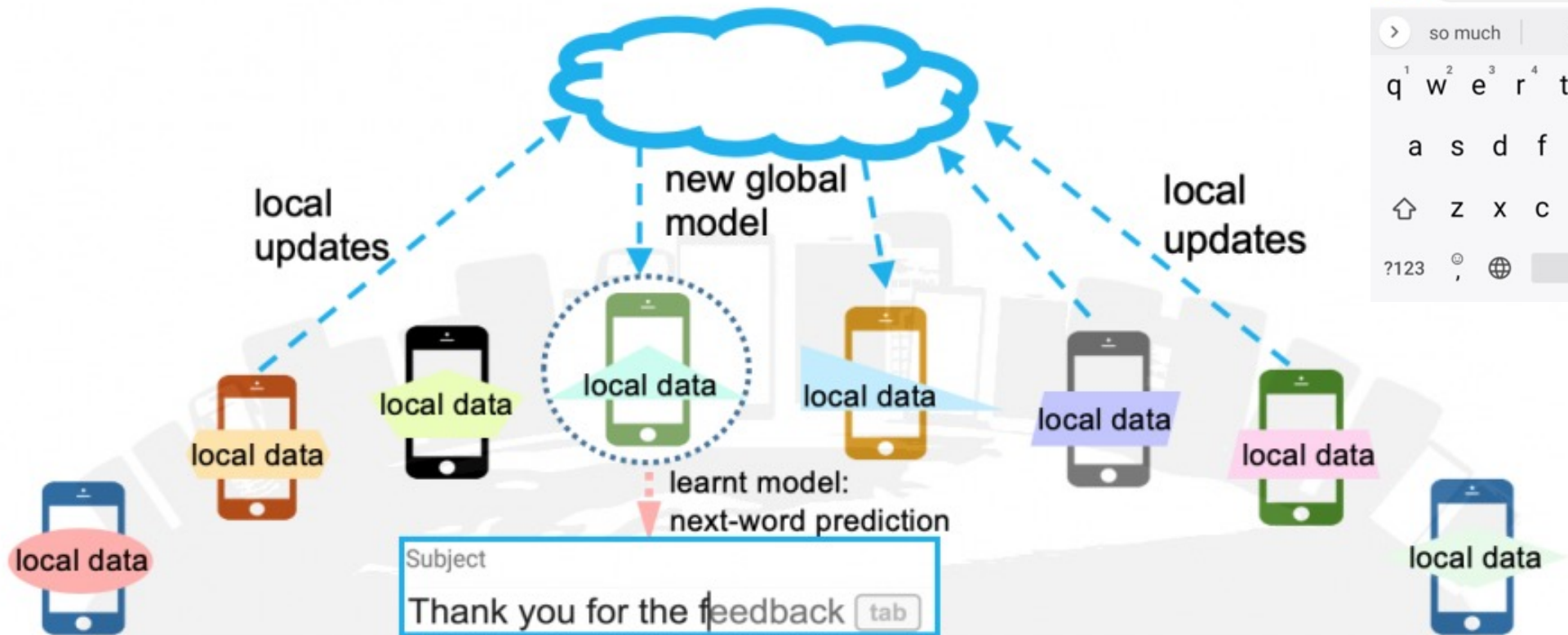
Federated learning works without the need to store user data in the cloud, but we're not stopping there. We've developed a Secure Aggregation protocol that uses cryptographic techniques so a coordinating server can only decrypt the average update if 100s or 1000s of users have participated — no individual phone's update can be inspected before averaging. It's the first protocol of its kind that is practical for deep-network-sized problems and real-world connectivity constraints. We designed Federated Averaging so the coordinating server only needs the average update, which allows Secure Aggregation to be used; however the protocol is general and can be applied to other problems as well. We're working hard on a production implementation of this protocol and expect to deploy it for Federated Learning applications in the near future.

Our work has only scratched the surface of what is possible. Federated Learning can't solve all machine learning problems (for example, learning to recognize different dog breeds by training on carefully labeled examples), and for many other models the necessary training data is already stored in the cloud (like training spam filters for Gmail). So Google will continue to advance the state-of-the-art for cloud-based ML, but we are also committed to ongoing research to expand the range of problems we can solve with Federated Learning. Beyond Gboard query suggestions, for example, we hope to improve the language models that power your keyboard based on what you actually type on your phone (which can have a style all its own) and photo rankings based on what kinds of photos people look at, share, or delete.

Applying Federated Learning requires machine learning practitioners to adopt new tools and a new way of thinking: model development, training, and evaluation with no direct access to or labeling of raw data, with communication cost as a limiting factor. We believe the user benefits of Federated Learning make tackling the technical challenges worthwhile, and are publishing our work with hopes of a widespread conversation within the machine learning community.

Jakub Konečný, H. Brendan McMahan, Felix X. Yu, Peter Richtárik, Ananda Theertha Suresh, Dave Bacon
**Federated Learning: Strategies for Improving Communication Efficiency**
*NIPS Private Multi-Party Machine Learning Workshop, 2016*

Jakub Konečný, H. Brendan McMahan, Daniel Ramage, Peter Richtárik
**Federated Optimization: Distributed Machine Learning for On-Device Intelligence**
*arXiv:1610.02527, 2016*

# The First Federated Learning App: Next-Word Prediction

**Federated Learning** is collaborative machine learning from private data stored across a (large) number of clients/devices (e.g., hospitals, phones)

# Peter Richtarik ✎

FOLLOWING

Professor, KAUST
Verified email at kaust.edu.sa - Homepage

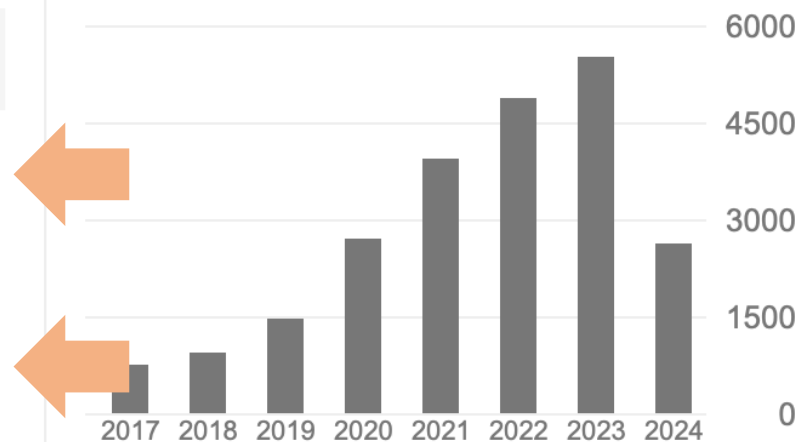optimization    machine learning    federated learning    deep learning
computer science

**Cited by**                                                    VIEW ALL

|            | All   | Since 2019 |
|------------|-------|------------|
| Citations  | 24873 | 21252      |
| h-index    | 69    | 63         |
| i10-index  | 180   | 171        |



2017 2018 2019 2020 2021 2022 2023 2024

**Public access**                                               VIEW ALL

0 articles                                                      43 articles

not available                                                   available

Based on funding mandates

| TITLE | CITED BY | YEAR |
|-------|----------|------|
| Federated learning: Strategies for improving communication efficiency <br> J Konecný, HB McMahan, FX Yu, P Richtárik, AT Suresh, D Bacon <br> arXiv preprint arXiv:1610.05492 8 | 3856 | 2016 |
| Federated learning: Strategies for improving communication efficiency <br> J Konečný, HB McMahan, FX Yu, P Richtárik, AT Suresh, D Bacon <br> arXiv preprint arXiv:1610.05492 | 2716 | 2016 |
| Federated optimization: Distributed machine learning for on-device intelligence <br> J Konečný, HB McMahan, D Ramage, P Richtárik <br> arXiv preprint arXiv:1610.02527 | 2091 | 2016 |
| Iteration complexity of randomized block-coordinate descent methods for minimizing a composite function <br> P Richtarik, M Takáč <br> Mathematical Programming 144 (2), 1-38 | 860 | 2014 |

# My Team: 100+ Papers on Federated Learning

## Peter Richtárik

Professor of Computer Science

KAUST

Address: Office 3145, Bldg 12, 4700 KAUST, Thuwal 23955-6900, Saudi Arabia
E-mail: peter.richtarik@kaust.edu.sa

All papers are listed below in reverse chronological order in which they appeared online.

## Prepared in 2024

[258] Kai Yi, Timur Kharisov, Igor Sokolov, and Peter Richtárik
**Cohort squeeze: Beyond a single communication round per cohort in cross-device federated learning**
Federated Learning Paper
[arXiv] [method: SPPM-AS]

[257] Georg Meinhardt, Kai Yi, Laurent Condat, and Peter Richtárik
**Prune at the clients, not the server: Accelerated sparse training in federated learning**
Federated Learning Paper
[arXiv] [method: Sparse-ProxSkip]

[256] Avetik Karagulyan, Egor Shulgin, Abdurakhmon Sadiev, and Peter Richtárik
**SPAM: Stochastic proximal point method with momentum variance reduction for non-convex cross-device federated learning**
Federated Learning Paper
[arXiv] [method: SPAM]

**Forbes**

AI

# The Next Generation Of Artificial Intelligence

**Rob Toews** Contributor ⓘ
*I write about the big picture of artificial intelligence.*

**Follow**

Oct 12, 2020, 09:22pm EDT

1. **Unsupervised Learning**
2. **Federated Learning**
3. **Transformers**
4. **Neural Network Compression**
5. **Generative AI**
6. **"System 2" Reasoning**

*NATIONAL ARTIFICIAL INTELLIGENCE*
*RESEARCH AND DEVELOPMENT*
*STRATEGIC PLAN*
*2023 UPDATE*

*A Report by the*

SELECT COMMITTEE ON ARTIFICIAL INTELLIGENCE
*of the*
NATIONAL SCIENCE AND TECHNOLOGY COUNCIL

May 2023

## Table of Contents

https://www.whitehouse.gov/wp-content/uploads/2023/05/National-Artificial-Intelligence-Research-and-Development-Strategic-Plan-2023-Update.pdf

# Part 2
# Introduction

# Optimization Problem

$$\min_{x \in \mathbb{R}^d} f(x) \stackrel{\text{def}}{=} \frac{1}{n} \sum_{i=1}^{n} f_i(x)$$

# parallel machines

# model parameters / features

Loss on local data $\mathcal{D}_i$ stored on machine $i$

$$f_i(x) := \mathbb{E}_{\xi \sim \mathcal{D}_i} \left[ f_i(x, \xi) \right]$$

**!** It takes $\tau_i$ seconds for worker $i$ to compute $\nabla f_i(x, \xi)$, where $\xi \sim \mathcal{D}_i$ $\qquad 0 < \tau_1 \leq \tau_2 \leq \cdots \leq \tau_n$

It takes $\theta_i$ seconds for worker $i$ to communicate $g \in \mathbb{R}^d$ to the server

Find a (possibly random) vector $\hat{x} \in \mathbb{R}^d$ such that $\mathbb{E}\left[ \|\nabla f(\hat{x})\|^2 \right] \leq \varepsilon$

# Parallel Computing Architecture

$x$ gets updated by the server

**Server**



$x$

$x$

$x$

$\nabla f_1(x, \xi)$

$\nabla f_2(x, \xi)$

$\nabla f_3(x, \xi)$

**Worker 1**

**Worker 2**

**Worker 3**

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

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

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

$\nabla f_1(x, \xi)$ compute time $= \tau_1$ secs

$\nabla f_2(x, \xi)$ compute time $= \tau_2$ secs

$\nabla f_3(x, \xi)$ compute time $= \tau_3$ secs

# Three Types of Heterogeneity

| | |
|---|---|
| **Data** | data distributions $\mathcal{D}_1, \ldots, \mathcal{D}_n$ can be different |
| **Compute** | compute times $\tau_1, \ldots, \tau_n$ are nonzero and can be different |
| **Communication** | communication times $\theta_1, \ldots, \theta_n$ are nonzero and can be different |

# Typical Assumptions

**1**  $\inf f \in \mathbb{R}$

**2**  $f_i(x) := \mathbb{E}_{\xi \sim \mathcal{D}_i} \left[ f_i(x, \xi) \right]$

Gradient of local functions is Lipschitz:

$$\max_{i \in \{1, \ldots, n\}} \sup_{x \neq y} \frac{\|\nabla f_i(x) - \nabla f_i(y)\|}{\|x - y\|} \leq L$$

Stochastic gradients have bounded variance:

$$\max_{i \in \{1, \ldots, n\}} \sup_{x \in \mathbb{R}^d} \mathbb{E}_{\xi \sim \mathcal{D}_i} \left[ \|\nabla f_i(x, \xi) - \mathbb{E}_{\xi \sim \mathcal{D}_i} \left[ \nabla f_i(x, \xi) \right] \|^2 \right] \leq \sigma^2$$

# Our Papers on Optimal Parallel SGD

**Optimal Time Complexities of Parallel Stochastic Optimization Methods Under a Fixed Computation Model**

5/2023

**Shadowheart SGD: Distributed Asynchronous SGD with Optimal Time Complexity Under Arbitrary Computation and Communication Heterogeneity**

2/2024

**Freya PAGE: First Optimal Time Complexity for Large-Scale Nonconvex Finite-Sum Optimization with Heterogeneous Asynchronous Computations**

5/2024

**On the Optimal Time Complexities in Decentralized Stochastic Asynchronous Optimization**

5/2024

# Our Papers

**5/2023**

Rennala SGD
Malenia SGD
Acc. Rennala SGD

Alexander Tyurin and P.R.
**Optimal time complexities of parallel stochastic optimization methods under a fixed computation model**
*NeurIPS 2023*

... *computation* (and/or data) **heterogeneity**

**2/2024**

Shadowheart SGD

Alexander Tyurin, Marta Pozzi, Ivan Ilin and P.R.
**Shadowheart SGD: Distributed asynchronous SGD with optimal time complexity under arbitrary computation and communication heterogeneity**
*arXiv:2402.04785, 2024*

... *communication* (and computation) **heterogeneity**

[Rennala SGD as a special case]

**5/2024**

Freya PAGE
Freya SGD

Alexander Tyurin, Kaja Gruntkowska, and P.R.
**Freya PAGE: First optimal time complexity for large-scale nonconvex finite-sum optimization with heterogeneous asynchronous computations**
*arXiv:2405.1554, 2024*

... *computation heterogeneity for* **finite-sum** *problems*

in the large-scale regime: $m \geq n^2$

**5/2024**

Fragile SGD, Amelie SGD + accelerated variants

Alexander Tyurin and P.R.
**On the optimal time complexities in decentralized stochastic asynchronous optimization**
*arXiv:2405.16218, 2024*

... *computation and communication heterogeneity in the* **decentralized setup**

# Peter, What About the Weird Algorithm Names?



Rennala, Queen of the Full Moon is a Legend Boss in Elden Ring. Though not a demigod, Rennala is one of the shardbearers who resides in the Academy of Raya Lucaria. Rennala is a powerful sorceress, head of the Carian Royal family, and erstwhile leader of the Academy.

# Optimal Parallel Stochastic Gradient Methods

| | Data Heterogeneity ($\mathcal{D}_i$ different) | Compute Heterogeneity ($\tau_i$ different) | Communication Heterogeneity ($\theta_i$ different) | Smooth Nonconvex | Smooth Convex | Infinite / Finite Sum? | Supports Decentralized Setup? | Optimal Time Complexity? |
|---|---|---|---|---|---|---|---|---|
| **Rennala SGD** <br> Tyurin & R (NeurIPS '23) | ✖ | ✔ | 0 | ✔ | | **Inf** | ✖ | ✔ |
| **Malenia SGD** <br> Tyurin & R (NeurIPS '23) | ✔ | ✔ | 0 | ✔ | | **Inf** | ✖ | ✔ |
| **Accelerated Rennala SGD** <br> Tyurin & R (NeurIPS '23) | ✖ | ✔ | 0 | | ✔ | **Inf** | ✖ | ✔ |
| **Shadowheart SGD** <br> Tyurin, Pozzi, Ilin & R '24 | ✖ | ✔ | ✔ | ✔ | | **Inf** | ✖ | ✔ |
| **Freya PAGE** <br> Tyurin, Gruntkowska & R '24 | ✖ | ✔ | 0 | ✔ | | **Finite** | ✖ | ✔ <br> big data regime |
| **Freya SGD** <br> Tyurin, Gruntkowska & R '24 | ✖ | ✔ | 0 | ✔ | | **Finite** | ✖ | ✖ |
| **Fragile SGD** <br> Tyurin & R '24 | ✖ | ✔ | ✔ | ✔ | | **Inf** | ✔ | nearly |
| **Amelie SGD** <br> Tyurin & R '24 | ✔ | ✔ | ✔ | ✔ | | **Inf** | ✔ | ✔ |

# Part 3
# Previous Approaches
# to Parallelizing SGD

# Hero SGD

Algorithmic idea: The fastest worker does it all!

# (Fair) Minibatch SGD

Algorithmic idea: Each worker does one job only!

# Asynchronous SGD

Algorithmic idea: All workers are slaves and useful

published in NIPS 2011

## NeurIPS 2020 Test of Time Award

# HOGWILD!: A Lock-Free Approach to Parallelizing Stochastic Gradient Descent

**Feng Niu**
leonn@cs.wisc.edu

**Benjamin Recht**
brecht@cs.wisc.edu

**Christopher Ré**
chrisre@cs.wisc.edu

**Stephen J. Wright**
swright@cs.wisc.edu
Computer Sciences Department
University of Wisconsin-Madison
Madison, WI 53706

### Abstract

Stochastic Gradient Descent (SGD) is a popular algorithm that can achieve state-of-the-art performance on a variety of machine learning tasks. Several researchers have recently proposed schemes to parallelize SGD, but all require performance-destroying memory locking and synchronization. This work aims to show using novel theoretical analysis, algorithms, and implementation that SGD can be implemented *without any locking*. We present an update scheme called HOGWILD! which allows processors access to shared memory with the possibility of overwriting each other's work. We show that when the associated optimization problem is *sparse*, meaning most gradient updates only modify small parts of the decision variable, then HOGWILD! achieves a nearly optimal rate of convergence. We demonstrate experimentally that HOGWILD! outperforms alternative schemes that use locking by an order of magnitude.

## 1 Introduction

With its small memory footprint, robustness against noise, and rapid learning rates, Stochastic Gradient Descent (SGD) has proved to be well suited to data-intensive machine learning tasks [3, 5, 24]. However, SGD's scalability is limited by its inherently sequential nature; it is difficult to parallelize. Nevertheless, the recent emergence of inexpensive multicore processors and mammoth, web-scale data sets has motivated researchers to develop several clever parallelization schemes fo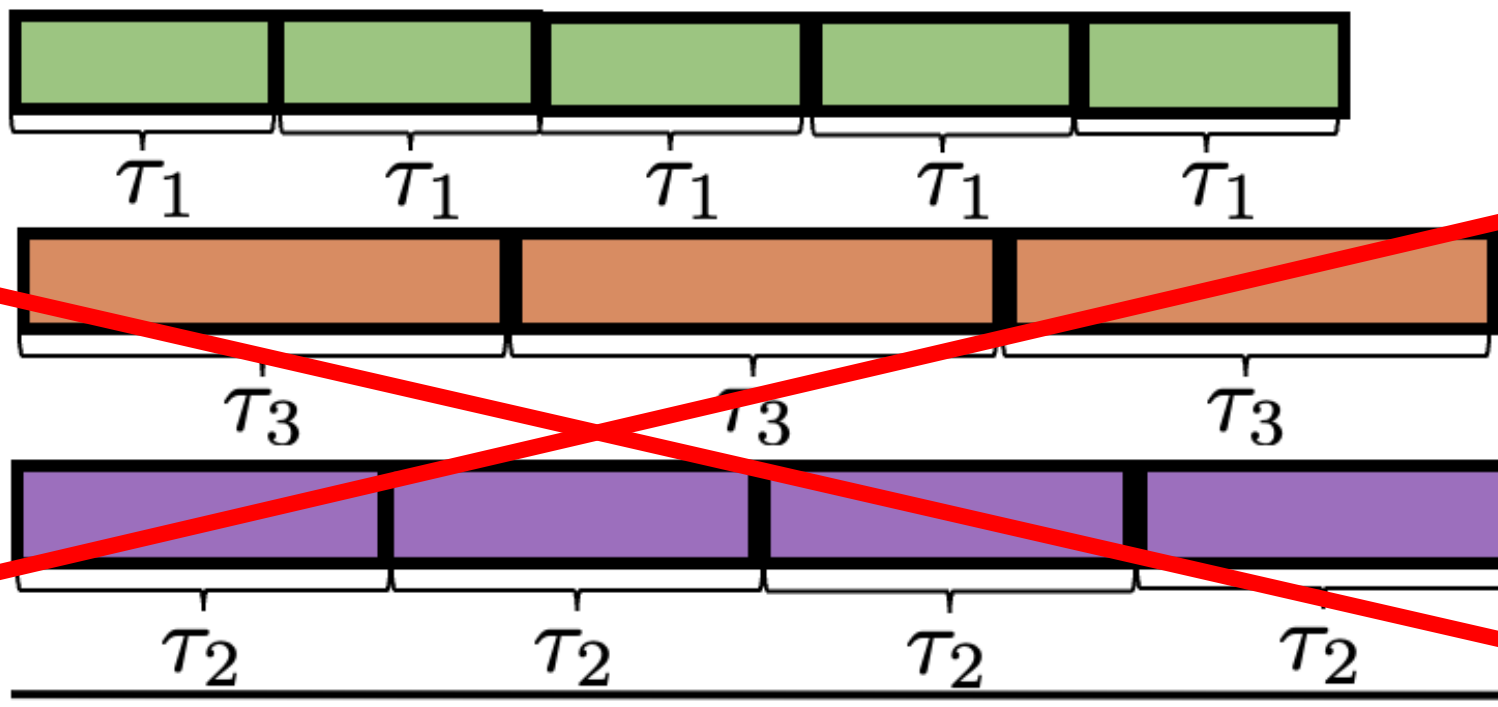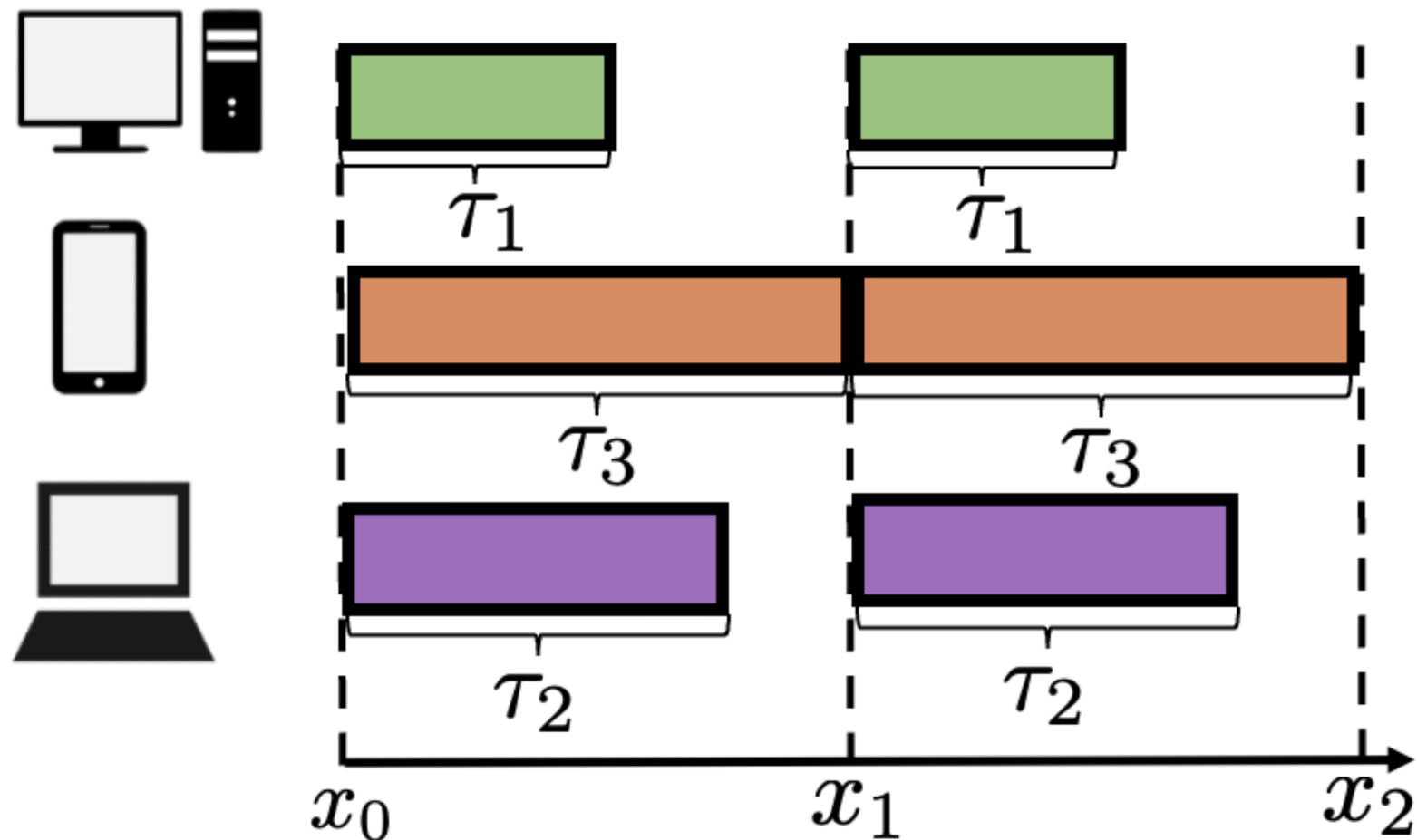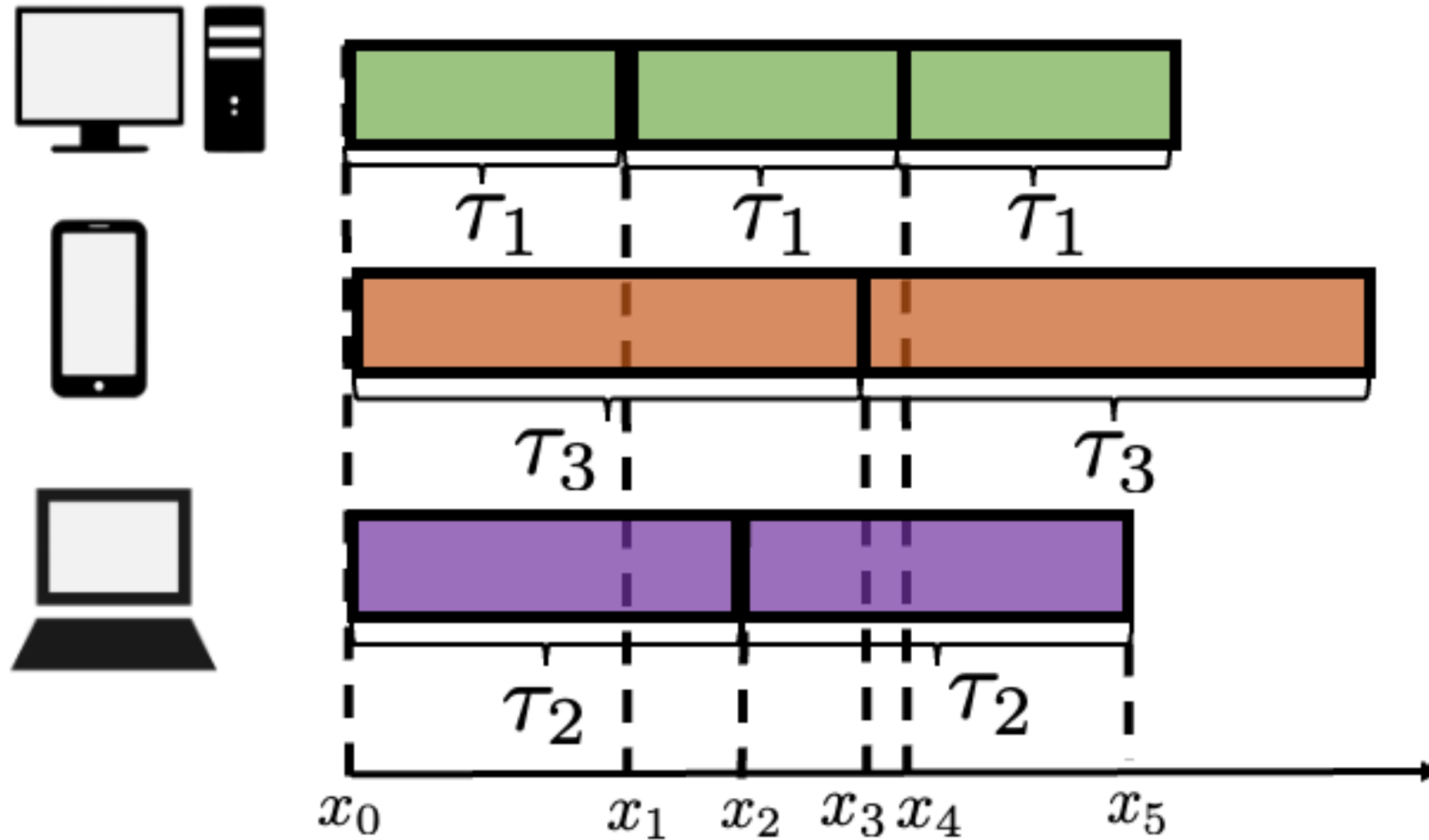r SGD [4, 10, 12, 16, 27]. As many large data sets are currently pre-processed in a MapReduce-like parallel-processing framework, much of the recent work on parallel SGD has focused naturally on MapReduce implementations. MapReduce is a powerful tool developed at Google for extracting information from huge logs (e.g., "find all the urls from a 100TB of Web data") that was designed to ensure fault tolerance and to simplify the maintenance and programming of large clusters of machines [9]. But MapReduce is not ideally suited for online, numerically intensive data analysis. Iterative computation is difficult to express in MapReduce, and the overhead to ensure fault tolerance can result in dismal throughput. Indeed, even Google researchers themselves suggest that other systems, for example Dremel, are more appropriate than MapReduce for data analysis tasks [20].

For some data sets, the sheer size of the data dictates that one use a cluster of machines. However, there are a host of problems in which, after appropriate preprocessing, the data necessary for statistical analysis may consist of a few terabytes or less. For such problems, one can use a single inexpensive work station as opposed to a hundred thousand dollar cluster. Multicore systems have significant performance advantages, including (1) low latency and high throughput shared main memory (a processor in such a system can write and read the shared physical memory at over 12GB/s with latency in the tens of nanoseconds); and (2) high bandwidth off multiple disks (a thousand-dollar RAID

1

# Our Inspiration: Two Beautiful Papers

## Asynchronous SGD Beats Minibatch SGD Under Arbitrary Delays

Konstantin Mishchenko    Francis Bach    Mathieu Even    Blake Woodworth

DI ENS, Ecole normale supérieure,
Université PSL, CNRS, INRIA
75005 Paris, France

### Abstract

The existing analysis of asynchronous stochastic gradient descent (SGD) degrades dramatically when any delay is large, giving the impression that performance depends primarily on the delay. On the contrary, we prove much better guarantees for the same asynchronous SGD algorithm regardless of the delays in the gradients, depending instead just on the number of parallel devices used to implement the algorithm. Our guarantees are strictly better than the existing analyses, and we also argue that asynchronous SGD outperforms synchronous minibatch SGD in the settings we consider. For our analysis, we introduce a novel recursion based on "virtual iterates" and delay-adaptive stepsizes, which allow us to derive state-of-the-art guarantees for both convex and non-convex objectives.

### 1 Introduction

We consider solving stochastic optimization problems of the form

$$\min_{\mathbf{x} \in \mathbb{R}^d} \{F(\mathbf{x}) := \mathbb{E}_{\xi \sim \mathcal{D}} f(\mathbf{x}; \xi)\}, \tag{1}$$

which includes machine learning (ML) training objectives, where $f(\mathbf{x}; \xi)$ represents the loss of a model parameterized by $\mathbf{x}$ on the datum $\xi$. Depending on the application, $\mathcal{D}$ could represent a finite dataset of size $n$ or a population distribution. In recent years, such stochastic optimization problems have continued to grow rapidly in size, both in terms of the dimension $d$ of the optimization variable—i.e., the number of model parameters in ML—and in terms of the quantity of data—i.e., the number of samples $\xi_1, \ldots, \xi_n \sim \mathcal{D}$ being used. With $d$ and $n$ regularly reaching the tens or hundreds of billions, it is increasingly necessary to use parallel optimization algorithms to handle the large scale and to benefit from data stored on different machines.

There are many ways of employing parallelism to solve (1), but the most popular approaches in practice are first-order methods based on stochastic gradient descent (SGD). At each iteration, SGD employs stochastic estimates of $\nabla F$ to update the parameters as $\mathbf{x}_k = \mathbf{x}_{k-1} - \gamma_k \nabla f(\mathbf{x}_{k-1}; \xi_{k-1})$ for an i.i.d. sample $\xi_{k-1} \sim \mathcal{D}$. Given $M$ machines capable of computing these stochastic gradient estimates $\nabla f(\mathbf{x}; \xi)$ in parallel, one approach to parallelizing SGD is what we call "Minibatch SGD." This refers to a synchronous, parallel algorithm that dispatches the current parameters $\mathbf{x}_{k-1}$ to each of the $M$ machines, waits while they compute and communicate back their gradient estimates $\mathbf{g}_{k-1}^1, \ldots, \mathbf{g}_{k-1}^M$, and then takes a minibatch SGD step $\mathbf{x}_k = \mathbf{x}_{k-1} - \gamma_k \cdot \frac{1}{M} \sum_{m=1}^M \mathbf{g}_{k-1}^m$. This is a natural idea with long history [16, 18, 55] and it is commonly used in practice [e.g., 22]. However, since Minibatch SGD waits for all $M$ of the machines to finish computing their gradient estimates before updating, it proceeds only at the speed of the *slowest* machine.

There are several possible sources of delays: nodes may have heterogeneous hardware with different computational throughputs [23, 25], network latency can slow the communication of gradients, and

---

## Sharper Convergence Guarantees for Asynchronous SGD for Distributed and Federated Learning

Anastasia Koloskova    Sebastian U. Stich    Martin Jaggi
EPFL              CISPA*            EPFL
anastasia.koloskova@epfl.ch  stich@cispa.de  martin.jaggi@epfl.ch

### Abstract

We study the asynchronous stochastic gradient descent algorithm for distributed training over $n$ workers which have varying computation and communication frequency over time. In this algorithm, workers compute stochastic gradients in parallel at their own pace and return those to the server without any synchronization. Existing convergence rates for this algorithm for non-convex smooth objectives depend on the maximum gradient delay $\tau_{\max}$ and show that an $\varepsilon$-stationary point is reached after $\mathcal{O}(\sigma^2 \varepsilon^{-2} + \tau_{\max} \varepsilon^{-1})$ iterations, where $\sigma$ denotes the variance of stochastic gradients.

In this work we obtain (i) a tighter convergence rate of $\mathcal{O}(\sigma^2 \varepsilon^{-2} + \sqrt{\tau_{\max} \tau_{avg}} \varepsilon^{-1})$ *without any change in the algorithm*, where $\tau_{avg}$ is the average delay, which can be significantly smaller than $\tau_{\max}$. We also provide (ii) a simple delay-adaptive learning rate scheme, under which asynchronous SGD achieves a convergence rate of $\mathcal{O}(\sigma^2 \varepsilon^{-2} + \tau_{avg} \varepsilon^{-1})$, and does not require any extra hyperparameter tuning nor extra communications. Our result allows to show *for the first time* that asynchronous SGD is *always faster* than mini-batch SGD. In addition, (iii) we consider the case of heterogeneous functions motivated by federated learning applications and improve the convergence rate by proving a weaker dependence on the maximum delay compared to prior works. In particular, we show that the heterogeneity term in convergence rate is only affected by the average delay within each worker.

### 1 Introduction

The stochastic gradient descent (SGD) algorithm [43, 13] and its variants (momentum SGD, Adam, etc.) form the foundation of modern machine learning and frequently achieve state of the art results. With recent growth in the size of models and available training data, parallel and distributed versions of SGD are becoming increasingly important [57, 17, 16]. Without those, modern state-of-the art language models [44], generative models [40, 41], and many others [50] would not be possible. In the distributed setting, also known as data-parallel training, optimization is distributed over many compute devices working in parallel (e.g. cores, or GPUs on a cluster) in order to speed up training. Every worker computes gradients on a subset of the training data, and the resulting gradients are aggregated (averaged) on a server.

The same type of SGD variants also form the core algorithms for federated learning applications [34, 24] where the training process is naturally distributed over many user devices, or clients, that keep their local data private, and only transfer (e.g. encrypted or differentially private) gradients to the server.

A rich literature exists on the convergence theory of above mentioned parallel SGD methods, see e.g. [17, 13] and references therein. Plain parallel SGD still faces many challenges in practice, motivat-

*CISPA Helmholtz Center for Information Security

---

arXiv: June 15, 2022                    arXiv: June 16, 2022

# Part 4
# Rennala SGD

Alexander Tyurin and P.R.
**Optimal time complexities of parallel stochastic optimization methods under a fixed computation model**
*NeurIPS 2023*

# Setup

## Optimal Parallel Stochastic Gradient Methods

| | Data Heterogeneity ($\mathcal{D}_i$ different) | Compute Heterogeneity ($\tau_i$ different) | Communication Heterogeneity ($\theta_i$ different) | Smooth Nonconvex | Smooth Convex | Infinite / Finite Sum? | Supports Decentralized Setup? | Optimal Time Complexity? |
|---|---|---|---|---|---|---|---|---|
| **Rennala SGD** Tyurin & R (NeurIPS '23) | ✗ | ✓ | 0 | ✓ | | Inf | ✗ | ✓ |
| **Malenia SGD** Tyurin & R (NeurIPS '23) | ✓ | ✓ | 0 | ✓ | | Inf | ✗ | ✓ |
| **Accelerated Rennala SGD** Tyurin & R (NeurIPS '23) | ✗ | ✓ | 0 | | ✓ | Inf | ✗ | ✓ |
| **Shadowheart SGD** Tyurin, Pozzi, Ilin & R '24 | ✗ | ✓ | ✓ | ✓ | | Inf | ✗ | ✓ |
| **Freya PAGE** Tyurin, Gruntkowska & R '24 | ✗ | ✓ | 0 | ✓ | | Finite | ✗ | ✓ big data regime |
| **Freya SGD** Tyurin, Gruntkowska & R '24 | ✗ | ✓ | 0 | ✓ | | Finite | ✗ | ✗ |
| **Fragile SGD** Tyurin & R '24 | ✗ | ✓ | ✓ | ✓ | | Inf | ✓ | nearly |
| **Amelie SGD** Tyurin & R '24 | ✓ | ✓ | ✓ | ✓ | | Inf | ✓ | ✓ |

# Rennala SGD

Algorithmic idea: Minibatch SGD with asynchronous minibatch collection

# Upper Bound

**Theorem (informal)**

Assume data homogeneity and zero communication times. Then Rennala SGD solves the problem in

Gradient of $f$ is $L$-Lipschitz

$\Delta := f(x^0) - \inf f$

Number of parallel machines

$$96 \times \min_{m \in \{1,\ldots,n\}} \left( \frac{1}{m} \sum_{i=1}^{m} \frac{1}{\tau_i} \right)^{-1} \left( \frac{L\Delta}{\varepsilon} + \frac{L\Delta\sigma^2}{\varepsilon^2 m} \right)$$

seconds.

Compute times

$0 < \tau_1 \leq \tau_2 \leq \cdots \leq \tau_n$

Algorithm outputs $\hat{x}$ such that $\mathbb{E}\left[\|\nabla f(\hat{x})\|^2\right] \leq \varepsilon$

$\sup_{x \in \mathbb{R}^d} \mathbb{E}_{\xi \sim \mathcal{D}}\left[\|\nabla f(x,\xi) - \nabla f(x)\|^2\right] \leq \sigma^2$

# Matching Lower Bound



**Upper Bound**

## Theorem (informal)

It is not possible to design a method that will find a solution faster than in

$$
\Omega \left( \min_{m \in \{1, \dots, n\}} \left( \frac{1}{m} \sum_{i=1}^{m} \frac{1}{\tau_i} \right)^{-1} \left( \frac{L\Delta}{\varepsilon} + \frac{L\Delta\sigma^2}{\varepsilon^2 m} \right) \right)
$$

seconds.

**Rennala SGD = first optimal parallel SGD**

# Classical Oracle: Keeps Track of # Iterations

Function class

Distribution governing noise

Oracle class

Algorithm class

**Protocol 1** Classical Oracle Protocol

1: **Input:** function $f \in \mathcal{F}$ oracle an

2: **for** $k = 0, \ldots, \infty$ **do**

3: $x^k = A^k(g^1, \ldots, g^k)$        $\triangleright$ $x^0 = A^0$ for $k = 0$.

4: $g^{k+1} = O(\ldots$

5: **end for**

stochastic gradient:
$$g^{k+1} = \nabla f(x^k, \xi^{k+1})$$

**Natural for sequential methods, where a single worker does all the work!**

**Iteration com**... measure):

$$\mathfrak{m}_{\mathrm{oracle}}\,(\mathcal{A}, \mathcal{F}) = \sup_{A \in \mathcal{A}} \sup_{f \in \mathcal{F}} \sup_{(O, \mathcal{D}) \in \mathcal{O}(f)} \inf \left\{ k \in \mathbb{N} \,\middle|\, \mathbb{E}\left[\|\nabla f(x^k)\|^2\right] \le \varepsilon \right\}$$

[Nemirovsky and Yudin, 1983]   [Nesterov, 2018]

[Carmon et al, 2020]   [Arjevani et al, 2022]

# New Oracle: Keeps Track of Time

**Protocol 2** Time Oracle Protocol

1: **Input:** functions $f \in \mathcal{F}$, oracle and distribution $(O, \mathcal{D}) \in \mathcal{O}$ [...] $A \in \mathcal{A}$
2: $s^0 = 0$
3: **for** $k = 0, \dots, \infty$ **do**
4:     $(t^{k+1}, x^k) = A^k(g^1, \dots, g^k),$                           $\triangleright t^{k+1} \geq t^k$
5:     $(s^{k+1}, g^{k+1}) = O(t^{k+1}, x^k, s^k$ [...]
6: **end for**

**Natural for parallel methods!**

**Iteration complex[...]**

$$\mathfrak{m}_{\mathrm{oracle}}(\mathcal{A}, \dots) = \dots \left\{ k \in \mathbb{N} \,\middle|\, \mathbb{E}\left[ \|\nabla f(x^k)\|^2 \right] \leq \varepsilon \right\}$$

$$\dots(f)$$

$$S_t := \left\{ k \in \mathbb{N} \cup \{0\} \,\middle|\, t^k \leq t \right\}$$

**Time complexity** [...] xity measure):

$$\mathfrak{m}_{\mathrm{time}}(\mathcal{A}, \mathcal{F}) := \inf_{A \in \mathcal{A}} \sup_{f \in \mathcal{F}} \sup_{(O, \mathcal{D}) \in \mathcal{O}(f)} \inf \left\{ t \geq 0 \,\middle|\, \mathbb{E}\left[ \inf_{k \in S_t} \|\nabla f(x^k)\|^2 \right] \leq \varepsilon \right\}$$

# Data Homogeneous Regime

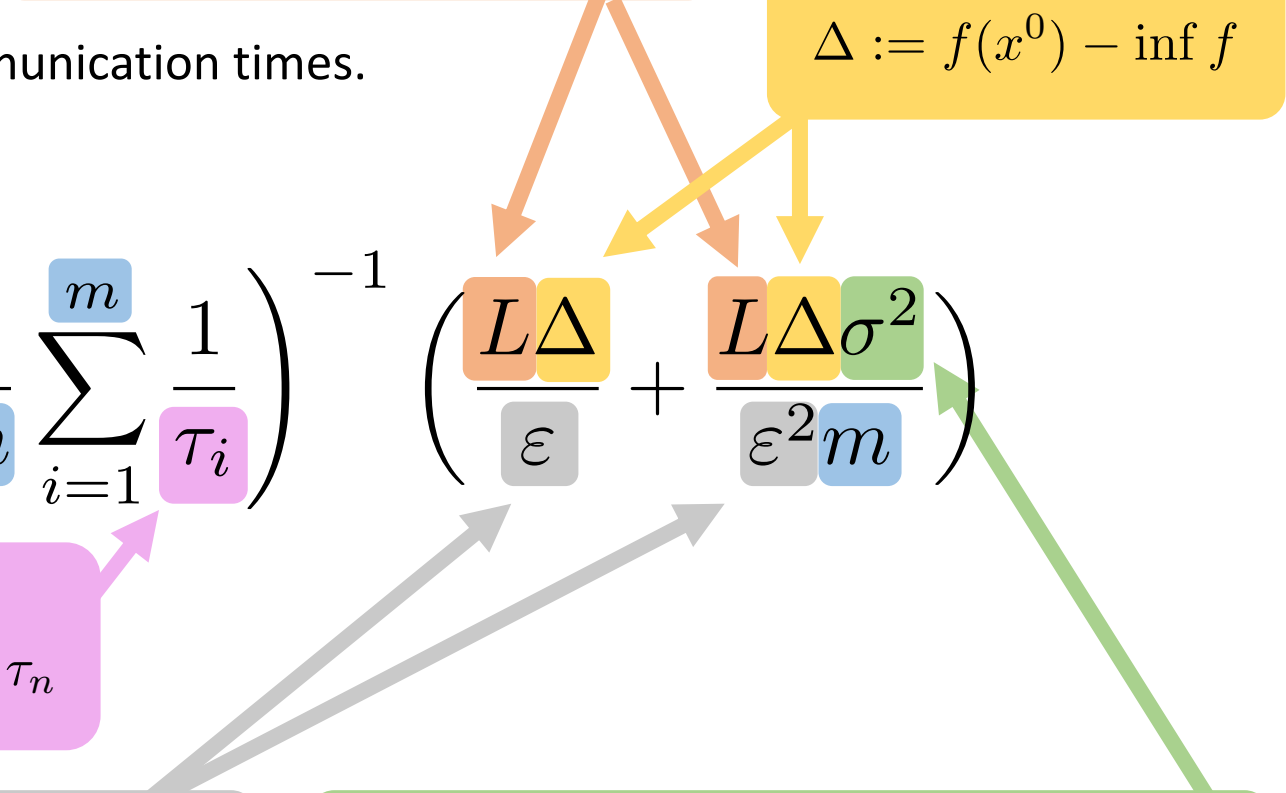| Method | Time Complexity |
| --- | --- |
| Minibatch SGD | $\tau_n \left( \frac{L\Delta}{\varepsilon} + \frac{\sigma^2 L\Delta}{n\varepsilon^2} \right)$ |
| Asynchronous SGD (Cohen et al., 2021) (Koloskova et al., 2022) (Mishchenko et al., 2022) | $\left( \frac{1}{n} \sum_{i=1}^{n} \frac{1}{\tau_i} \right)^{-1} \left( \frac{L\Delta}{\varepsilon} + \frac{\sigma^2 L\Delta}{n\varepsilon^2} \right)$ |
| Rennala SGD (Theorem 7.5) | $\min_{m \in [n]} \left[ \left( \frac{1}{m} \sum_{i=1}^{m} \frac{1}{\tau_i} \right)^{-1} \left( \frac{L\Delta}{\varepsilon} + \frac{\sigma^2 L\Delta}{m\varepsilon^2} \right) \right]$ |
| Lower Bound (Theorem 6.4) | $\min_{m \in [n]} \left[ \left( \frac{1}{m} \sum_{i=1}^{m} \frac{1}{\tau_i} \right)^{-1} \left( \frac{L\Delta}{\varepsilon} + \frac{\sigma^2 L\Delta}{m\varepsilon^2} \right) \right]$ |

# Experimental Results (Sample)

$$\tau_i = \sqrt{i} \text{ seconds}$$



Figure 3: # of workers $n = 10000$.
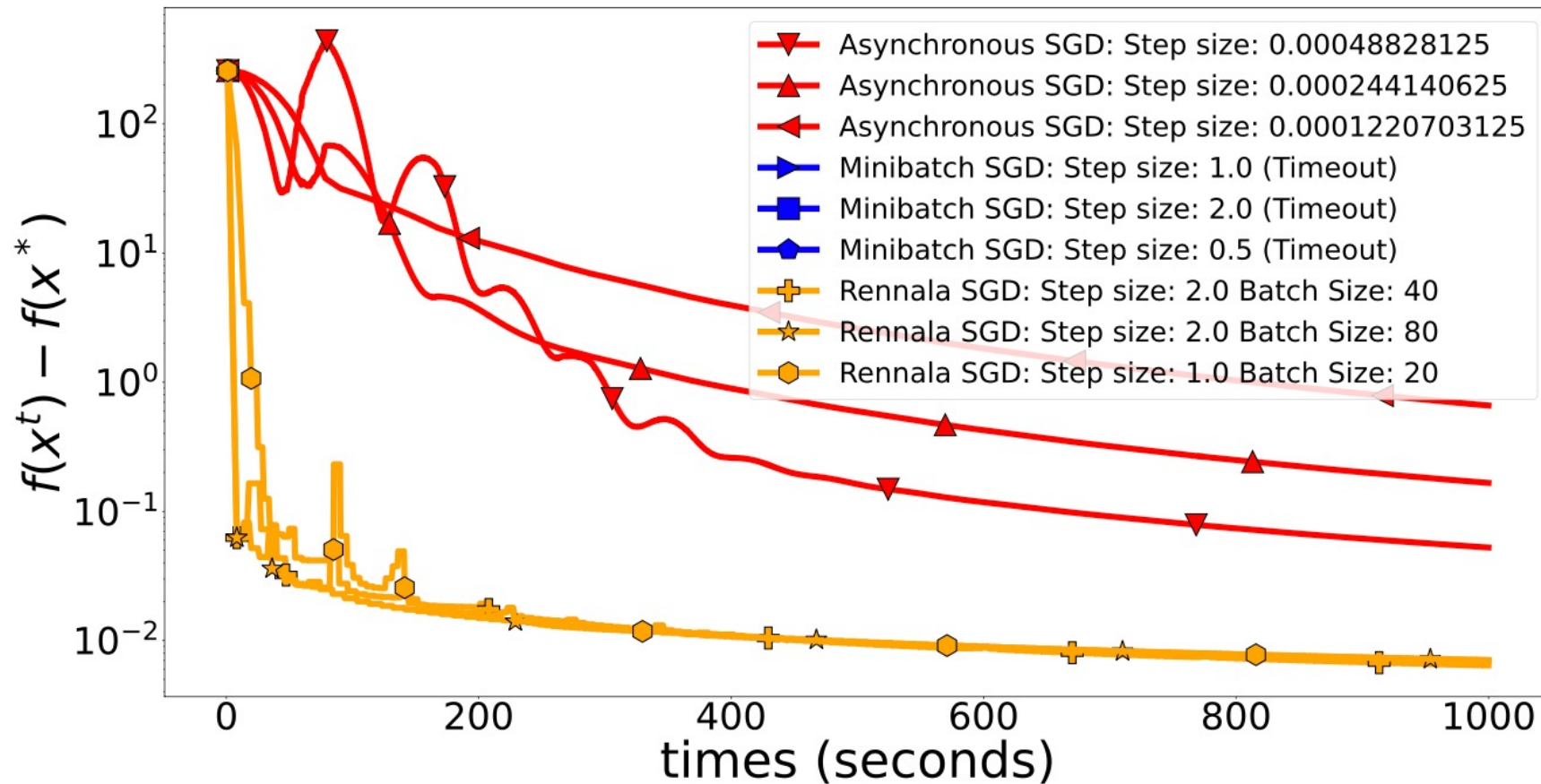
# Part 5
# Two Extensions

Alexander Tyurin and P.R.

Optimal time complexities of parallel stochastic optimization methods under a fixed computation model

*NeurIPS 2023*

# Extension 1
# Handling Data Heterogeneity
# (Malenia SGD)

# Malenia SGD: Setup

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

$$f_i(x) := \mathbb{E}_{\xi \sim \mathcal{D}_i} \left[ f_i(x, \xi) \right]$$

## Optimal Parallel Stochastic Gradient Methods

| | Data Heterogeneity ($\mathcal{D}_i$ different) | Compute Heterogeneity ($\tau_i$ different) | Communication Heterogeneity ($\theta_i$ different) | Smooth Nonconvex | Smooth Convex | Infinite / Finite Sum? | Supports Decentralized Setup? | Optimal Time Complexity? |
|---|---|---|---|---|---|---|---|---|
| **Rennala SGD** Tyurin & R (NeurIPS '23) | ✗ | ✓ | 0 | ✓ | | Inf | ✗ | ✓ |
| **Malenia SGD** Tyurin & R (NeurIPS '23) | ✓ | ✓ | 0 | ✓ | | Inf | ✗ | ✓ |
| **Accelerated Rennala SGD** Tyurin & R (NeurIPS '23) | ✗ | ✓ | 0 | | ✓ | Inf | ✗ | ✓ |
| **Shadowheart SGD** Tyurin, Pozzi, Ilin & R '24 | ✗ | ✓ | ✓ | ✓ | | Inf | ✗ | ✓ |
| **Freya PAGE** Tyurin, Gruntkowska & R '24 | ✗ | ✓ | 0 | ✓ | | Finite | ✗ | ✓ big data regime |
| **Freya SGD** Tyurin, Gruntkowska & R '24 | ✗ | ✓ | 0 | ✓ | | Finite | ✗ | ✗ |
| **Fragile SGD** Tyurin & R '24 | ✗ | ✓ | ✓ | ✓ | | Inf | ✓ | nearly |
| **Amelie SGD** Tyurin & R '24 | ✓ | ✓ | ✓ | ✓ | | Inf | ✓ | ✓ |

The distributions $\mathcal{D}_1, \ldots, \mathcal{D}_n$ are allowed to be different

# Malenia SGD

**Method 6** Malenia SGD

1: **Input:** starting point $x^0$, stepsize $\gamma$, parameter $S$
2: Run Method 7 in all workers
3: **for** $k = 0, 1, \ldots, K - 1$ **do**
4:      Init $g_i^k = 0$ and $B_i = 0$
5:      **while** $\left( \frac{1}{n} \sum_{i=1}^n \frac{1}{B_i} \right)^{-1} < \frac{S}{n}$ **do**
6:          Wait for the next worker
7:          Receive gradient, iteration index, worker's index $(g, k', i)$
8:          **if** $k' = k$ **then**
9:              $g_i^k = g_i^k + g$
10:            $B_i = B_i + 1$
11:          **end if**
12:          Send $(x^k, k)$ to the worker
13:      **end while**
14:      $g^k = \frac{1}{n} \sum_{i=1}^n \frac{1}{B_i} g_i^k$
15:      $x^{k+1} = x^k - \gamma g^k$
16: **end for**

Minibatch size

$$S = \max \left\{ \left\lceil \frac{\sigma^2}{\varepsilon} \right\rceil, n \right\}$$

**Method 7** Worker's Infinite Loop

1: Init $g = 0$, $k' = -1$, and worker's index $i$
2: **while** True **do**
3:      Send $(g, k', i)$ to the server
4:      Receive $(x^k, k)$ from the server
5:      $k' = k$
6:      $g = \widehat{\nabla} f_i(x^k; \xi), \quad \xi \sim \mathcal{D}$
7: **end while**

# (Nonconvex) Data Heterogeneous Regime

| Method | Time Complexity |
|---|---|
| Minibatch SGD | $\tau_n \left( \frac{L\Delta}{\varepsilon} + \frac{\sigma^2 L\Delta}{n\varepsilon^2} \right)$ |
| Malenia SGD (Theorem A.4) | $\tau_n \frac{L\Delta}{\varepsilon} + \left( \frac{1}{n} \sum_{i=1}^{n} \tau_i \right) \frac{\sigma^2 L\Delta}{n\varepsilon^2}$ |
| Lower Bound (Theorem A.2) | $\tau_n \frac{L\Delta}{\varepsilon} + \left( \frac{1}{n} \sum_{i=1}^{n} \tau_i \right) \frac{\sigma^2 L\Delta}{n\varepsilon^2}$ |

# Extension 2
# Handling the Convex Regime
# (Accelerated Rennala SGD)

# Accelerated Rennala SGD: Setup

## Optimal Parallel Stochastic Gradient Methods

| | Data Heterogeneity ($\mathcal{D}_i$ different) | Compute Heterogeneity ($\tau_i$ different) | Communication Heterogeneity ($\theta_i$ different) | Smooth Nonconvex | Smooth Convex | Infinite / Finite Sum? | Supports Decentralized Setup? | Optimal Time Complexity? |
|---|---|---|---|---|---|---|---|---|
| **Rennala SGD** Tyurin & R (NeurIPS '23) | ✘ | ✔ | 0 | ✔ | | Inf | ✘ | ✔ |
| **Malenia SGD** Tyurin & R (NeurIPS '23) | ✔ | ✔ | 0 | ✔ | | Inf | ✘ | ✔ |
| **Accelerated Rennala SGD** Tyurin & R (NeurIPS '23) | ✘ | ✔ | 0 | | ✔ | Inf | ✘ | ✔ |
| **Shadowheart SGD** Tyurin, Pozzi, Ilin & R '24 | ✘ | ✔ | ✔ | ✔ | | Inf | ✘ | ✔ |
| **Freya PAGE** Tyurin, Gruntkowska & R '24 | ✘ | ✔ | 0 | ✔ | | Finite | ✘ | ✔ big data regime |
| **Freya SGD** Tyurin, Gruntkowska & R '24 | ✘ | ✔ | 0 | ✔ | | Finite | ✘ | ✘ |
| **Fragile SGD** Tyurin & R '24 | ✘ | ✔ | ✔ | ✔ | | Inf | ✔ | nearly |
| **Amelie SGD** Tyurin & R '24 | ✔ | ✔ | ✔ | ✔ | | Inf | ✔ | ✔ |

# Convex (Data Homogeneous) Regime

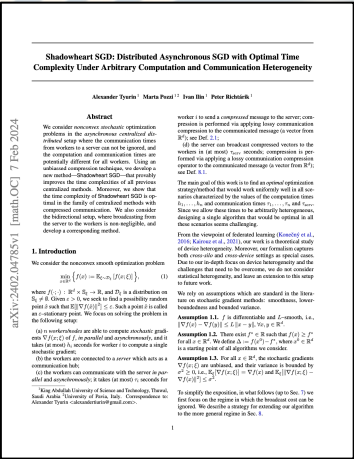| Method | Time Complexity |
|---|---|
| Minibatch SGD | $\tau_n \left( \min \left\{ \frac{\sqrt{L}R}{\sqrt{\varepsilon}}, \frac{M^2 R^2}{\varepsilon^2} \right\} + \frac{\sigma^2 R^2}{n\varepsilon^2} \right)$ |
| Asynchronous SGD (Mishchenko et al., 2022) | $\left( \frac{1}{n} \sum_{i=1}^{n} \frac{1}{\tau_i} \right)^{-1} \left( \frac{LR^2}{\varepsilon} + \frac{\sigma^2 R^2}{n\varepsilon^2} \right)$ |
| (Accelerated) Rennala SGD (Theorems B.9 and B.11) | $\min_{m \in [n]} \left[ \left( \frac{1}{m} \sum_{i=1}^{m} \frac{1}{\tau_i} \right)^{-1} \left( \min \left\{ \frac{\sqrt{L}R}{\sqrt{\varepsilon}}, \frac{M^2 R^2}{\varepsilon^2} \right\} + \frac{\sigma^2 R^2}{m\varepsilon^2} \right) \right]$ |
| Lower Bound (Theorem B.4) | $\min_{m \in [n]} \left[ \left( \frac{1}{m} \sum_{i=1}^{m} \frac{1}{\tau_i} \right)^{-1} \left( \min \left\{ \frac{\sqrt{L}R}{\sqrt{\varepsilon}}, \frac{M^2 R^2}{\varepsilon^2} \right\} + \frac{\sigma^2 R^2}{m\varepsilon^2} \right) \right]$ |
| Lower Bound (Section M) (Woodworth et al., 2018) | $\tau_1 \min \left\{ \frac{\sqrt{L}R}{\sqrt{\varepsilon}}, \frac{M^2 R^2}{\varepsilon^2} \right\} + \left( \frac{1}{n} \sum_{i=1}^{n} \frac{1}{\tau_i} \right)^{-1} \frac{\sigma^2 R^2}{n\varepsilon^2}$ |

$\nabla f$ is $L$-Lipschitz, $f$ is $M$-Lipschitz, and $\|x^0 - x^\star\| \le R$

# The End

# Further Extensions

# Shadowheart SGD

## Optimal Parallel SGD
## under Compute Heterogeneity
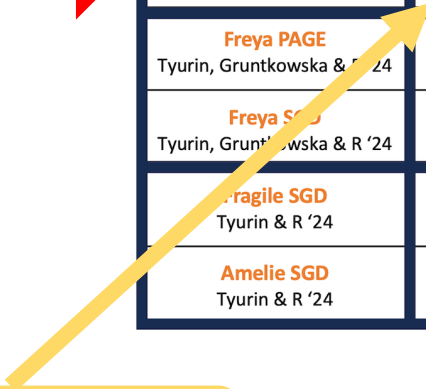## & Communication Heterogeneity

# Shadowheart SGD: Setup

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

$$f_i(x) := \mathbb{E}_{\xi \sim \mathcal{D}_i}[f_i(x, \xi)]$$

## Optimal Parallel Stochastic Gradient Methods

| | Data Heterogeneity ($\mathcal{D}_i$ different) | Compute Heterogeneity ($\tau_i$ different) | Communication Heterogeneity ($\theta_i$ different) | Smooth Nonconvex | Smooth Convex | Infinite / Finite Sum? | Supports Decentralized Setup? | Optimal Time Complexity? |
|---|---|---|---|---|---|---|---|---|
| **Rennala SGD** Tyurin & R (NeurIPS '23) | ✗ | ✓ | 0 | ✓ | | Inf | ✗ | ✓ |
| **Malenia SGD** Tyurin & R (NeurIPS '23) | ✓ | ✓ | 0 | ✓ | | Inf | ✗ | ✓ |
| **Accelerated Rennala SGD** Tyurin & R (NeurIPS '23) | ✗ | ✓ | 0 | | ✓ | Inf | ✗ | ✓ |
| **Shadowheart SGD** Tyurin, Pozzi, Ilin & R '24 | ✗ | ✓ | ✓ | ✓ | | Inf | ✗ | ✓ |
| **Freya PAGE** Tyurin, Gruntkowska & '24 | ✗ | ✓ | 0 | ✓ | | Finite | ✗ | ✓ big data regime |
| **Freya SGD** Tyurin, Gruntkowska & R '24 | ✗ | ✓ | 0 | ✓ | | Finite | ✗ | ✗ |
| **Fragile SGD** Tyurin & R '24 | ✗ | ✓ | ✓ | ✓ | | Inf | ✓ | nearly |
| **Amelie SGD** Tyurin & R '24 | ✓ | ✓ | ✓ | ✓ | | Inf | ✓ | ✓ |

$\mathcal{D}_1 = \cdots = \mathcal{D}_n$

Communication costs $\theta_1, \ldots, \theta_n$ are nonzero (and possibly different)

# Shadowheart SGD

Table 1: **Time Complexities of Centralized Distributed Algorithms.** Assume that it takes at most $h_i$ seconds to worker $i$ to calculate a stochastic gradient and $\dot{\tau}_i$ seconds to send *one coordinate/float* to server. Abbreviations: $L$ = smoothness constant, $\varepsilon$ = error tolerance, $\Delta = f(x^0) - f^*$, $n$ = # of workers, $d$ = dimension of the problem. We take the Rand$K$ compressor with $K = 1$ (Def. C.1) (as an example) in QSGD and Shadowheart SGD. Due to Property 5.2, the choice $K = 1$ is optimal for Shadowheart SGD up to a constant factor.
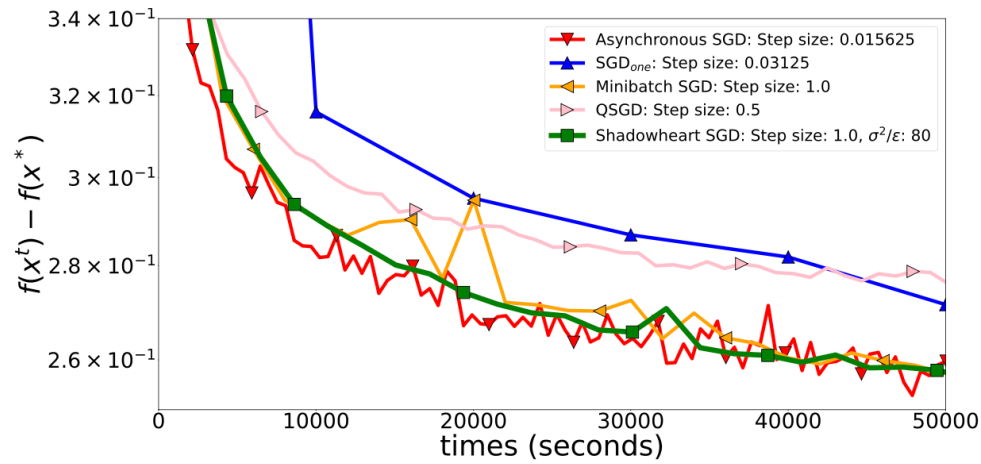
| Method | Time Complexity | $\max\{h_n, \dot{\tau}_n\} \to \infty,$ $\max\{h_i, \dot{\tau}_i\} < \infty \, \forall i < n$ (the last worker is slow) | Time Complexities in Some Regimes $h_i = h, \dot{\tau}_i = \dot{\tau} \, \forall i \in [n]$ (equal performance) | Numerical Comparison[b] $\sigma^2/\varepsilon =$ | | |
|---|---|---|---|---|---|---|
| | | | | 1 | $10^3$ | $10^6$ |
| Minibatch SGD (see (3)) | $\max\limits_{i \in [n]} \max\{h_i, d\dot{\tau}_i\} \left( \frac{L\Delta}{\varepsilon} + \frac{\sigma^2 L\Delta}{n\varepsilon^2} \right)$ | $\infty$ (non-robust) | $\max\{h, d\dot{\tau}, \frac{d\dot{\tau}\sigma^2}{n\varepsilon}, \frac{h\sigma^2}{n\varepsilon}\} \frac{L\Delta}{\varepsilon}$ (worse, e.g., when $\dot{\tau}, d$ or $n$ large) | $\times 10^3$ | $\times 10^3$ | $\times 10^4$ |
| QSGD (see (7)) (Alistarh et al., 2017) (Khaled & Richtárik, 2020) | $\max\limits_{i \in [n]} \max\{h_i, \dot{\tau}_i\} \left( \left(\frac{d}{n} + 1\right) \frac{L\Delta}{\varepsilon} + \frac{d\sigma^2 L\Delta}{n\varepsilon^2} \right)$ | $\infty$ (non-robust) | $\geq \frac{dh\sigma^2}{n\varepsilon} \frac{L\Delta}{\varepsilon}$ (worse, e.g., when $\varepsilon$ small) | $\times 3$ | $\times 10^2$ | $\times 10^4$ |
| Rennala SGD (Tyurin & Richtárik, 2023c), Asynchronous SGD (e.g., (Mishchenko et al., 2022)) | $\geq \min\limits_{j \in [n]} \max\left\{ h_{\bar{\pi}_j}, d\dot{\tau}_{\bar{\pi}_j}, \frac{\sigma^2}{\varepsilon} \left( \sum\limits_{i=1}^{j} \frac{1}{h_{\bar{\pi}_i}} \right)^{-1} \right\} \frac{L\Delta}{\varepsilon}$ [a] | $< \infty$ (robust) | $\geq \max\left\{ h, d\dot{\tau}, \frac{h\sigma^2}{n\varepsilon} \right\} \frac{L\Delta}{\varepsilon}$ (worse, e.g., when $\dot{\tau}, d$ or $n$ large) | $\times 10^2$ | $\times 10$ | $\times 1.5$ |
| Shadowheart SGD (see (9) and Alg. 1) (Corollary 4.4) | $t^*(d - 1, \sigma^2/\varepsilon, [h_i, \dot{\tau}_i]_1^n) \frac{L\Delta}{\varepsilon}$ [c] | $< \infty$ (robust) | $\max\left\{ h, \dot{\tau}, \frac{d\dot{\tau}}{n}, \sqrt{\frac{d\dot{\tau}h\sigma^2}{n\varepsilon}}, \frac{h\sigma^2}{n\varepsilon} \right\} \frac{L\Delta}{\varepsilon}$ | $\times 1$ | $\times 1$ | $\times 1$ |

The time complexity of Shadowheart SGD is not worse than the time complexity of the competing centralized methods (see Sec. 6), and is *strictly* better in many regimes. We show that (12) is the *optimal time complexity* in the family of centralized methods with compression (see Sec. 7).
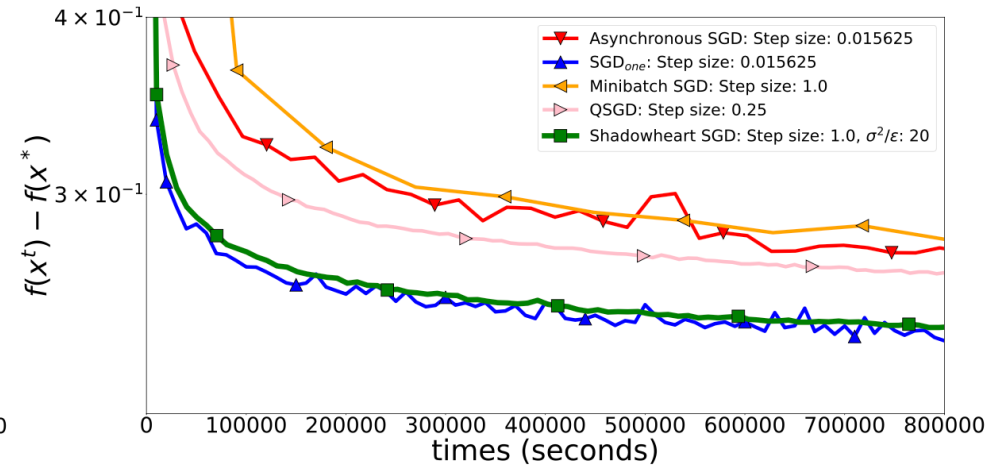
[a] Upper bound time complexities are not derived for Rennala SGD and Asynchronous SGD. However, we can derive the lower bound using Theorem N.5 with $\omega = 0$. One should take $d\dot{\tau}_i$ instead of $\tau_i$ when apply Theorem N.5 because these methods send $d$ coordinates. $\bar{\pi}$ is a permutation that sorts $\max\{h_i, d\dot{\tau}_i\} : \max\{h_{\bar{\pi}_1}, d\dot{\tau}_{\bar{\pi}_1}\} \leq \cdots \leq \max\{h_{\bar{\pi}_n}, d\dot{\tau}_{\bar{\pi}_n}\}$

[b] We numerically compute time complexities for $d = 10^6$, $n = 10^3$, $h_i \sim U(0.1, 1)$, $\dot{\tau}_i \sim U(0.1, 1)$ (uniform i.i.d.), and three noise regimes $\sigma^2/\varepsilon \in \{1, 10^3, 10^6\}$. We report the factors by which the time complexities of the competing methods are worse compared to the time complexity of our method Shadowheart SGD. So, for example, Minibatch SGD, QSGD and Asynchronous SGD can be worse by the factors $\times 10^4$, $\times 10^4$, and $\times 10^2$, respectively.
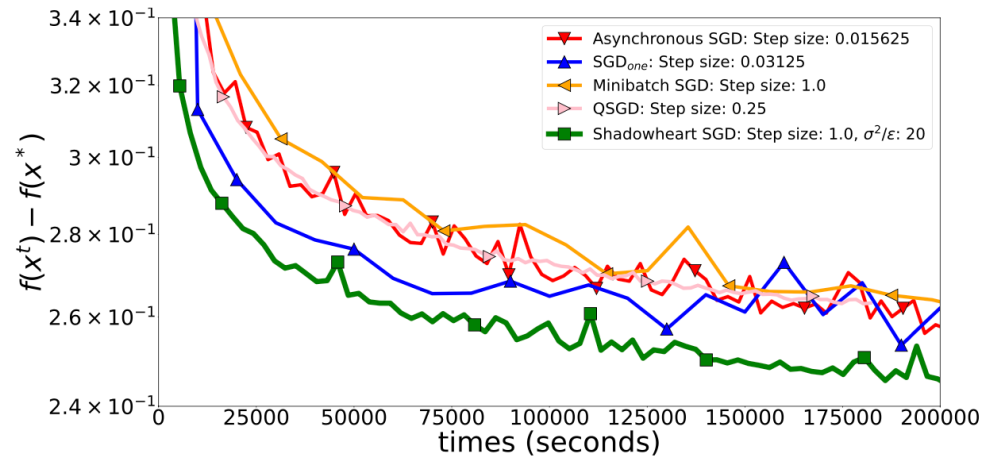
[c] The mapping $t^*$ is defined in Def. 4.2.

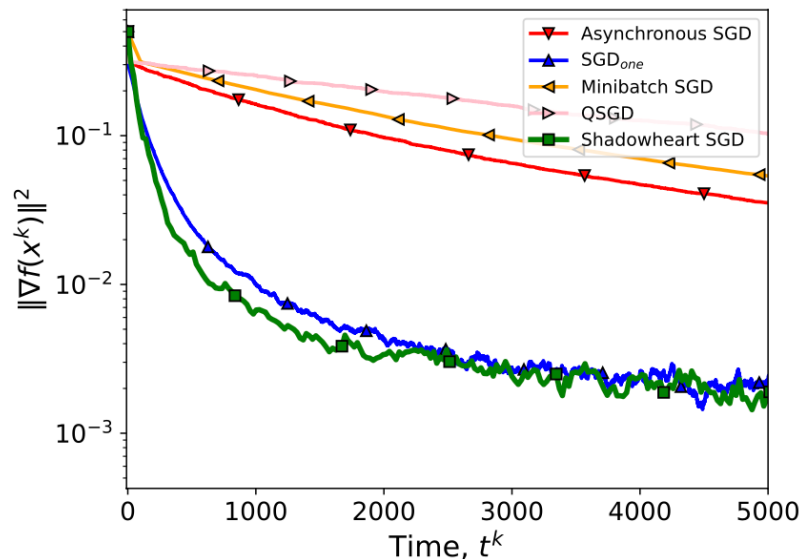**Fast** communication: $\dot{\theta}_i = \frac{\sqrt{i}}{d}$

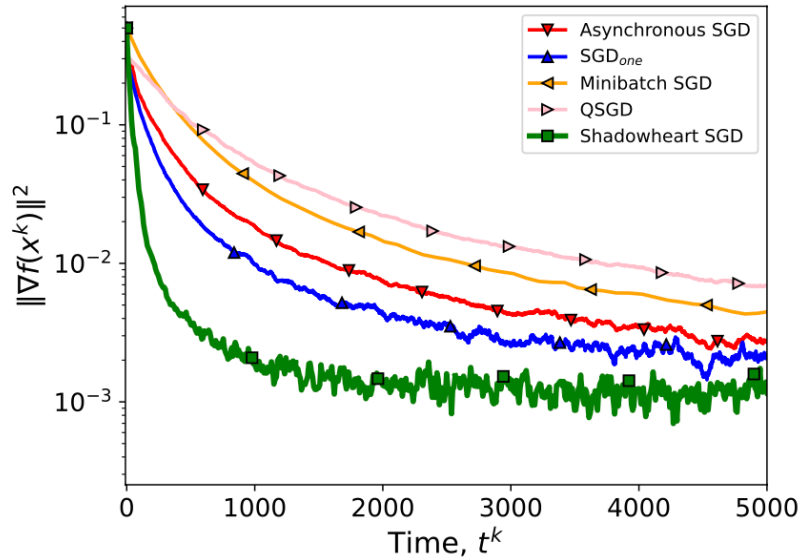**Slow** communication: $\dot{\theta}_i = \frac{\sqrt{i}}{d^{1/2}}$

**Medium-speed** communication: $\dot{\theta}_i = \frac{\sqrt{i}}{d^{3/4}}$

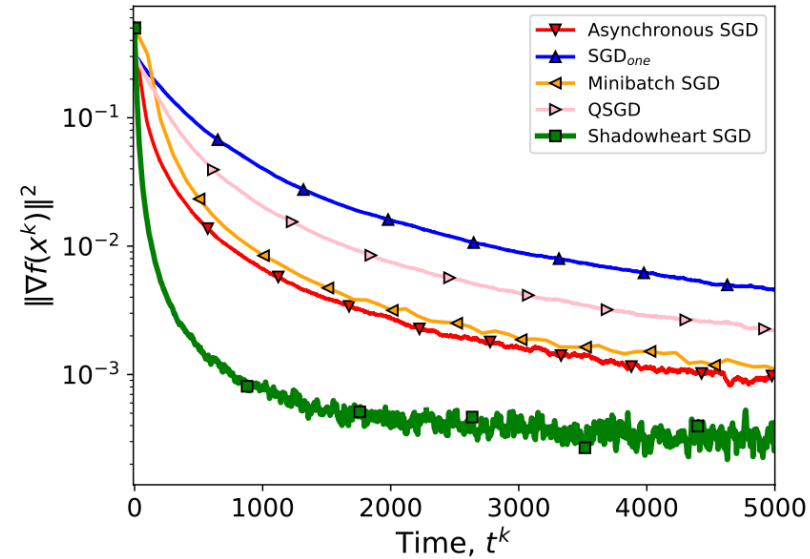Computation times: $\tau_i = \sqrt{i}$ for all machines $i = 1, \ldots, n$

# Shadowheart SGD: Adding More Workers...



(a) $n = 10$      (b) $n = 10^2$      (c) $n = 10^3$

$$\tau_i^k, \dot{\theta}_i^k \sim \text{Uniform}(0.1, 1) \text{ for all } i \in \{1, \ldots, n\} \text{ and } k \geq 0$$

Shadowheart

Freya

# Freya PAGE

## Optimal Parallel SGD
## for Large-Scale Finite-Sum Problems

# Freya PAGE: Setup

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

$$f_i(x) := \mathbb{E}_{\xi \sim \mathcal{D}_i} \left[ f_i(x, \xi) \right]$$

## Optimal Parallel Stochastic Gradient Methods

| | Data Heterogeneity ($\mathcal{D}_i$ different) | Compute Heterogeneity ($\tau_i$ different) | Communication Heterogeneity ($\theta_i$ different) | Smooth Nonconvex | Smooth Convex | Infinite / Finite Sum? | Supports Decentralized Setup? | Optimal Time Complexity? |
|---|---|---|---|---|---|---|---|---|
| **Rennala SGD** <br> Tyurin & R (NeurIPS '23) | ✗ | ✓ | 0 | ✓ | | Inf | ✗ | ✓ |
| **Malenia SGD** <br> Tyurin & R (NeurIPS '23) | ✓ | ✓ | 0 | ✓ | | Inf | ✗ | ✓ |
| **Accelerated Rennala SGD** <br> Tyurin & R (NeurIPS '23) | ✗ | ✓ | 0 | | ✓ | Inf | ✗ | ✓ |
| **Shadowheart SGD** <br> Tyurin, Pozzi, Ilin & R '24 | ✗ | ✓ | ✓ | ✓ | | Inf | ✗ | ✓ |
| **Freya PAGE** <br> Tyurin, Gruntkowska & R '24 | ✗ | ✓ | 0 | ✓ | | **Finite** | ✗ | ✓ <br> big data regime |
| **Freya SGD** <br> Tyurin, Gruntkowska & R '24 | ✗ | ✓ | 0 | ✓ | | Finite | ✗ | ✗ |
| **Fragile SGD** <br> Tyurin & R '24 | ✗ | ✓ | ✓ | ✓ | | Inf | ✓ | nearly |
| **Amelie SGD** <br> Tyurin & R '24 | ✓ | ✓ | ✓ | ✓ | | Inf | ✓ | ✓ |

$$\mathcal{D}_1 = \cdots = \mathcal{D}_n$$

$$\mathcal{D}_i = \text{uniform distribution over } m \text{ outcomes}$$

# PAGE: Optimal Serial SGD for Finite-Sum Nonconvex Optimization

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

$$f_i(x) := \mathbb{E}_{\xi \sim \mathcal{D}_i} \left[ f_i(x, \xi) \right]$$

$$\mathcal{D}_1 = \cdots = \mathcal{D}_n$$

$$\mathcal{D}_i = \text{uniform distribution over } m \text{ outcomes}$$

$$\min_{x \in \mathbb{R}^d} \left\{ f(x) = \frac{1}{m} \sum_{i=1}^{m} f_i(x) \right\}$$

(after butchering/redefining notation)

Zhize Li, Hongyan Bao, Xiangliang Zhang, and P.R.
**PAGE: A simple and optimal probabilistic gradient estimator for nonconvex optimization**
*ICML 2021*

Table 1: Comparison of the *worst-case time complexity* guarantees of methods that work with asynchronous computations in the setup from Section 1 (up to smoothness constants). We assume that $\tau_i \in [0, \infty]$ is the bound on the times required to calculate one stochastic gradient $\nabla f_j$ by worker $i$, $\tau_1 \leq \ldots \leq \tau_n$, and $m \geq n \log n$. Abbr: $\delta^0 := f(x^0) - f^*$, $m = \#$ of data samples, $n = \#$ of workers, $\varepsilon = $ error tolerance.

| Method | Worst-Case Time Complexity | Comment |
|---|---|---|
| Hero GD (Soviet GD) | $\tau_1 m \frac{\delta^0}{\varepsilon} \quad \left(\tau_n \frac{m}{n} \frac{\delta^0}{\varepsilon}\right)$ | Suboptimal |
| Hero PAGE (Soviet PAGE) [Li et al., 2021] | $\tau_1 m + \tau_1 \frac{\delta^0}{\varepsilon}\sqrt{m} \quad \left(\tau_n \frac{m}{n} + \tau_n \frac{\delta^0}{\varepsilon} \frac{\sqrt{m}}{n}\right)$ | Suboptimal |
| SYNTHESIS [Liu et al., 2022] | — | Limitations: bounded gradient assumption, calculates the full gradients[a], suboptimal.[b] |
| Asynchronous SGD [Koloskova et al., 2022] [Mishchenko et al., 2022] | $\frac{\delta^0}{\varepsilon}\left(\left(\sum_{i=1}^{n}\frac{1}{\tau_i}\right)^{-1}\left(\frac{\sigma^2}{\varepsilon}+n\right)\right)$ | Limitations: $\sigma^2$–bounded variance assumption, suboptimal when $\varepsilon$ is small. |
| Rennala SGD [Tyurin and Richtárik, 2023] | $\frac{\delta^0}{\varepsilon}\min_{j\in[n]}\left(\left(\sum_{i=1}^{j}\frac{1}{\tau_i}\right)^{-1}\left(\frac{\sigma^2}{\varepsilon}+j\right)\right)$ | Limitations: $\sigma^2$–bounded variance assumption, suboptimal when $\varepsilon$ is small. |
| Freya PAGE (Theorems 7 and 8) | $\min_{j\in[n]}\left(\left(\sum_{i=1}^{j}\frac{1}{\tau_i}\right)^{-1}(m+j)\right)$ $+\frac{\delta^0}{\varepsilon}\min_{j\in[n]}\left(\left(\sum_{i=1}^{j}\frac{1}{\tau_i}\right)^{-1}(\sqrt{m}+j)\right)^{(c)}$ | Optimal in the large-scale regime, i.e., $\sqrt{m} \geq n$ (see Section 5) |
| Lower bound (Theorem 10) | $\min_{j\in[n]}\left(\left(\sum_{i=1}^{j}\frac{1}{\tau_i}\right)^{-1}(m+j)\right)$ $+\frac{\delta^0}{\sqrt{m}\varepsilon}\min_{j\in[n]}\left(\left(\sum_{i=1}^{j}\frac{1}{\tau_i}\right)^{-1}(m+j)\right)$ | — |

Freya PAGE has *universally* better guarantees than all previous methods: the dependence on $\varepsilon$ is $\mathcal{O}(1/\varepsilon)$ (unlike Rennala SGD and Asynchronous SGD), the dependence on $\{\tau_i\}$ is harmonic-like and robust to slow workers (robust to $\tau_n \to \infty$) (unlike Soviet PAGE and SYNTHESIS), the assumptions are weak, and the time complexity of Freya PAGE is optimal when $\sqrt{m} \geq n$.

[a] In Line 3 of their Algorithm 3, they calculate the full gradient, assuming that it can be done for free and not explaining how.
[b] Their convergence rates in Theorems 1 and 3 depend on a bound on the delays $\Delta$, which in turn depends on the performance of the slowest worker. Our method does not depend on the slowest worker if it is too slow (see Section 4.3), which is required for optimality.
[c] We prove better time complexity in Theorem 6, but this result requires the knowledge of $\{\tau_i\}$ in advance, unlike Theorems 7 and 8.

**Algorithm 1** Freya PAGE

1: **Parameters:** starting point $x^0 \in \mathbb{R}^d$, learning rate $\gamma > 0$, minibatch size $S \in \mathbb{N}$, probability $p \in (0, 1]$, initialization $g^0 = \nabla f(x^0)$ using ComputeGradient($x^0$)    (Alg. 2)

2: **for** $k = 0, 1, \dots, K-1$ **do**

3:      $x^{k+1} = x^k - \gamma g^k$

4:      Sample $c^k \sim$ Bernoulli($p$)

5:      **if** $c^k = 1$ **then**                                        (with probability $p$)

6:          $\nabla f(x^{k+1}) = $ ComputeGradient($x^{k+1}$)                    (Alg. 2)

7:          $g^{k+1} = \nabla f(x^{k+1})$

8:      **else**                                               (with probability $1 - p$)

9:          $\frac{1}{S} \sum_{i \in \mathcal{S}^k} \left( \nabla f_i(x^{k+1}) - \nabla f_i(x^k) \right) = $ ComputeBatchDifference($S, x^{k+1}, x^k$)     (Alg. 3)

10:          $g^{k+1} = g^k + \frac{1}{S} \sum_{i \in \mathcal{S}^k} \left( \nabla f_i(x^{k+1}) - \nabla f_i(x^k) \right)$

11:      **end if**

12: **end for**

     (note): $\mathcal{S}^k$ is a set of i.i.d. indices that are sampled from $[m]$, *uniformly with replacement*, $\left| \mathcal{S}^k \right| = S$
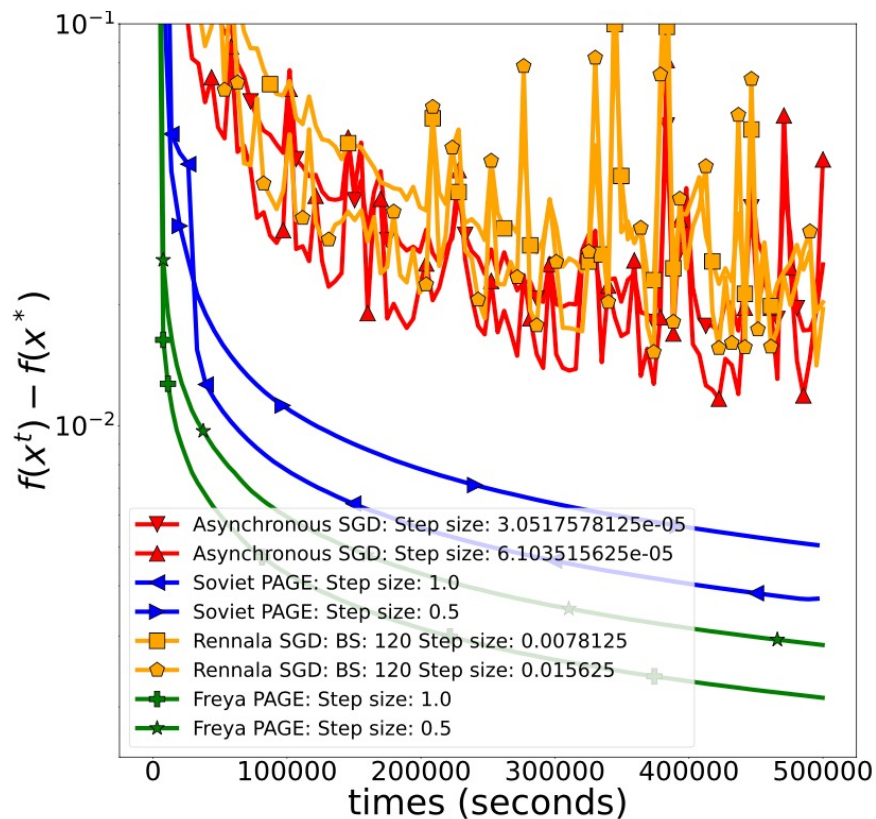
**Algorithm 2** ComputeGradient($x$)

1: **Input:** point $x \in \mathbb{R}^d$
2: Init $g = 0 \in \mathbb{R}^d$, set $\mathcal{M} = \emptyset$
3: Broadcast $x$ to all workers
4: For each worker $i \in [n]$, sample $j$ from $[m]$ uniformly and ask it to calculate $\nabla f_j(x)$
5: **while** $\mathcal{M} \neq [m]$ **do**
6:     Wait for $\nabla f_p(x)$ from a worker
7:     **if** $p \in [m] \backslash \mathcal{M}$ **then**
8:         $g \leftarrow g + \frac{1}{m}\nabla f_p(x)$
9:         Update $\mathcal{M} \leftarrow \mathcal{M} \cup \{p\}$
10:     **end if**
11:     Sample $j$ from $[m] \backslash \mathcal{M}$ uniformly and ask this worker to calculate $\nabla f_j(x)$
12: **end while**
13: Return $g = \frac{1}{m}\sum_{i=1}^{m}\nabla f_i(x)$

---

**Algorithm 3** ComputeBatchDifference($S, x, y$)

1: **Input:** batch size $S \in \mathbb{N}$, points $x, y \in \mathbb{R}^d$
2: Init $g = 0 \in \mathbb{R}^d$
3: Broadcast $x, y$ to all workers
4: For each worker, sample $j$ from $[m]$ uniformly and ask it to calculate $\nabla f_j(x) - \nabla f_j(y)$
5: **for** $i = 1, 2, \ldots, S$ **do**
6:     Wait for $\nabla f_p(x) - \nabla f_p(y)$ from a worker
7:     $g \leftarrow g + \frac{1}{S}(\nabla f_p(x) - \nabla f_p(y))$
8:     Sample $j$ from $[m]$ uniformly and ask this worker to calculate $\nabla f_j(x) - \nabla f_j(y)$
9: **end for**
10: Return $g$

Notes: i) the workers can aggregate $\nabla f_p$ locally, and the algorithm can call AllReduce once to collect all calculated gradients. ii) By splitting $[m]$ into blocks, instead of one $\nabla f_p$, we can ask the workers to calculate the sum of one block in Alg. 2 (and use a similar idea in Alg. 3).
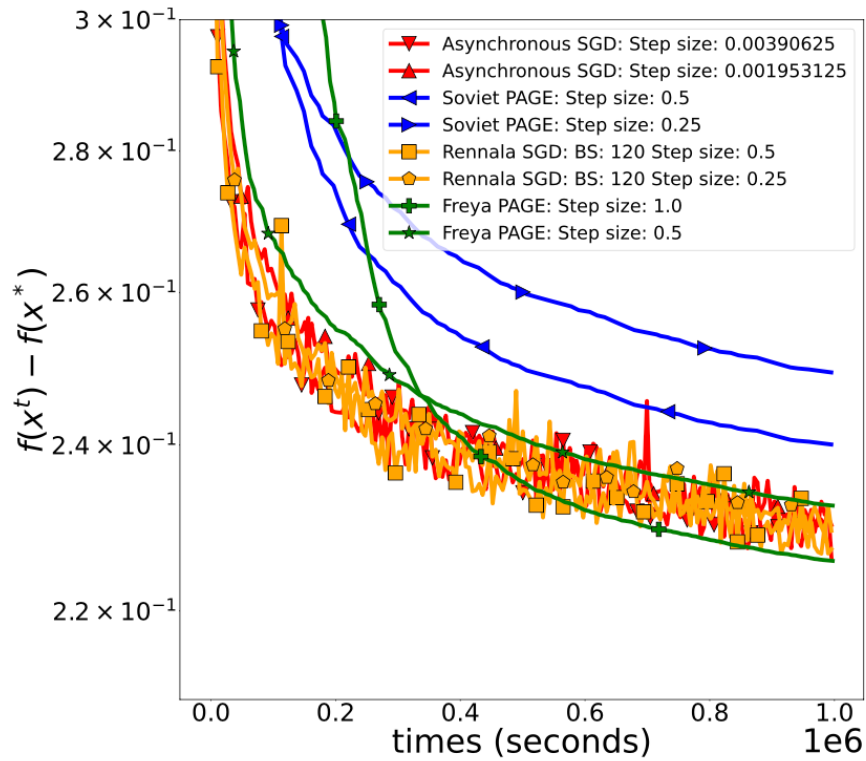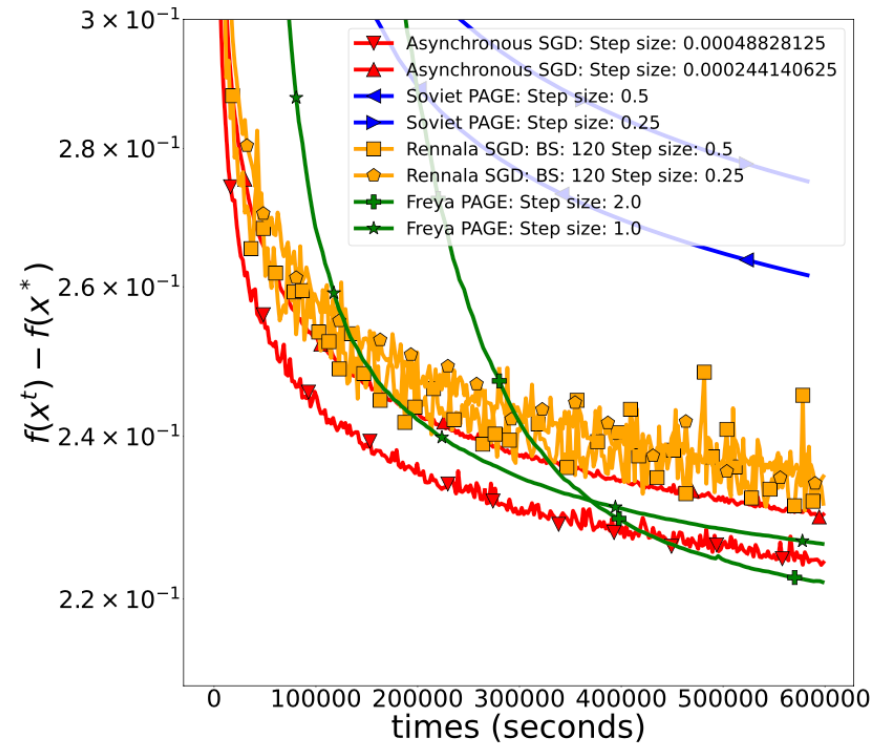
# Freya PAGE: Experiment 1



(a) $n = 1000$

(b) $n = 10000$

Figure 1: Experiments with nonconvex quadratic optimization tasks. We plot function suboptimality against elapsed time.

# Freya PAGE: Experiment 2



Figure 2: Experiments with the logistic regression problem on the MNIST dataset.

# Freya PAGE: Experiment 2

Table 2: Mean and variance of algorithm accuracies on the `MNIST` test set during the final 100K seconds of the experiments from Figure 2b.

| Method | Accuracy | Variance of Accuracy |
|---|---|---|
| Asynchronous SGD [Koloskova et al., 2022] [Mishchenko et al., 2022] | 92.60 | 5.85e-07 |
| Soviet PAGE [Li et al., 2021] | 92.31 | 1.62e-07 |
| Rennala SGD [Tyurin and Richtárik, 2023] | 92.37 | 3.12e-06 |
| **Freya PAGE** | **92.66** | **1.01e-07** |

# Amelie SGD

## Optimal Decentralized SGD
## under Computation & Communication Heterogeneity

# Decentralized Setup: Amelie SGD

| Method | The Worst-Case Time Complexity Guarantees | Comment |
|---|---|---|
| Minibatch SGD | $\frac{L\Delta}{\varepsilon} \max \left\{ \left(1 + \frac{\sigma^2}{n\varepsilon}\right) \max\{ \max\limits_{i,j\in[n]} \tau_{i\to j}, \max\limits_{i\in[n]} h_i \} \right\}$ | suboptimal if $\sigma^2/\varepsilon$ is large |
| RelaySGD, Gradient Tracking (Vogels et al., 2021) (Liu et al., 2024) | $\geq \frac{\max\limits_{i\in[n]} L_i \Delta}{\varepsilon} \frac{\sigma^2}{n\varepsilon} \max\limits_{i\in[n]} h_i$ | requires local $L_i$-smooth. of $f_i$, suboptimal if $\sigma^2/\varepsilon$ is large (even if $\max_{i\in[n]} L_i = L$) |
| Asynchronous SGD (Even et al., 2024) | — | requires similarity of the functions $\{f_i\}$, requires local $L_i$-smooth. of $f_i$ |
| Amelie SGD and Lower Bound (Thm. 7 and Cor. 2) | $\frac{L\Delta}{\varepsilon} \max \left\{ \max\limits_{i,j\in[n]} \tau_{i\to j}, \max\limits_{i\in[n]} h_i, \frac{\sigma^2}{n\varepsilon} \left( \frac{1}{n} \sum\limits_{i=1}^{n} h_i \right) \right\}$ | Optimal up to a constant factor |

# The End
# (for real)