





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



October 10-12 / Innopolis, Russia

Optimization & Machine Learning Lab @ KAUST



Part 1 Federated Learning





H Brendan McMahan



Jakub Konečný





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



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.

Q	umar	ni bur	ger n	nenu				12
Q	uman	ni bur	ger					ĸ
G	uma	ami b	urge					
~		۲				\$		۴
q	w	3 3	r	t ⁵ y ⁶	u ⁷	i	09	p°

To make Federated Learning possible, a had to overcome many algorithmic and technical challenges. In a typical machine learning, "stem, an optimization algorithm like Stochastic Gradient Descent (SGO) runs on a large dataset part, need homogeneously across servers in the cloud. Such highly iterative algorithms require low-latency, anthroughput connections to the training data. But in the Federated Learning setting, the data is districtly and across millions of devices in a highly uneven fashion. In addition, these devices have significant, ther-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 highdimensional sparsemut, a models which excel on problems like click-through-rate prediction.

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

Keith Bonawitz et al

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



The system then needs to communicate and aggregate a 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 stor, 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 Averaning applications in the near future.

Our work has only scratched the surface of what is possible. Federated Learning cart 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-theart 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



P	and the second second	Peter Richtarik 🖌		FOLLOWING	Cited by		VIEW ALL
	xe	Professor, <u>KAUST</u> Verified email at kaust edu sa - Homenage				All	Since 2019
	to	optimization machine learning federated learning deep learning	arning computer science	e	Citations h-index i10-index	24939 69 186	21428 64 177
	TITLE 🕒	:	CITED E	BY YEAR			6000
	Federated learr J Konečný, HB Mcl arXiv preprint arXiv	ning: Strategies for improving communication efficiency Mahan, FX Yu, P Richtárik, AT Suresh, D Bacon :1610.05492	55	03 2016		лI	4500
	Federated optin J Konečný, HB Mcl arXiv preprint arXiv	nization: Distributed machine learning for on-device intellige Mahan, D Ramage, P Richtárik :1610.02527	ence 22	67 2016	ad	Ш	- 1500
	Iteration comple composite funct P Richtarik, M Taká Mathematical Prog	exity of randomized block-coordinate descent methods for m tion ^{ič} ramming 144 (2), 1-38	ninimizing a 8	63 2014	2017 2018 2019 2020	0 2021 2022 2023	3 2024 U
	Generalized po M Journee, Y Neste Journal of Machine	wer method for sparse principal component analysis erov, P Richtárik, R Sepulchre Learning Research 11, 517-553	7:	51 2010	0 articles		43 articles
	Parallel coordin P Richtárik, M Taká Mathematical Prog	ate descent methods for big data optimization ^{ič} ramming 156 (1), 433-484	5	48 2016	Based on funding m	andates	avaliable
	Scaling distribut A Sapio, M Canini, 18th USENIX Symp	ted machine learning with {In-Network} aggregation CY Ho, J Nelson, P Kalnis, C Kim, A Krishnamurthy, posium on Networked Systems Design and Implementation (NSDI …	4	52 2021	Co-authors		EDIT
	Tighter theory fo A Khaled, K Mishch The 23rd Internatio	or local SGD on identical and heterogeneous data nenko, P Richtárik nal Conference on Artificial Intelligence and Statistics	4	52 2020	Martin Taká Mohamed b	č in Zayed Univers	sity o >
	SGD: General A RM Gower, N Loizo ICML 2019	Analysis and Improved Rates ou, X Qian, A Sailanbayev, E Shulgin, P Richtarik	4.	44 2019	Jakub Kone Research S Konstantin I Samsung	čný cientist, Google Mishchenko	>

My Team: 100+ Papers on Federated Learning



Peter Richtárik

Professor of Computer Science

King Abdullah University of Science and Technology (KAUST)

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



News Old News Papers Talks Video Talks Events Code Team Join Bio Teaching Consulting

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

Prepared in 2024

[261] Artavazd Maranjyan, Omar Shaikh Omar, and Peter Richtárik MindFlayer: Efficient asynchronous parallel SGD in the presence of heterogeneous and random worker compute times Asynchronous Optimization [arXiv] [method: MindFlayer SGD, Vecna SGD]

[260] Hanmin Li and Peter Richtárik On the convergence of FedProx with extrapolation and inexact prox Federated Learning Paper [arXiv] [method: FedExProx]

[259] Eduard Gorbunov, Nazarii Tupitsa, Sayantan Choudhury, Alen Aliev, Peter Richtárik, Samuel Horváth, and Martin Takáč Methods for convex \$(L_0,L_1)\$-smooth optimization: clipping, acceleration, and adaptivity [arXiv] [method: L0L1-GD, L0L1-GD-PS, L0L1-STM, L0L1-AdGD, L0L1-SGD, L0L1-SGD-PS]

[258] Kai Yi, Timur Kharisov, Igor Sokolov, and Peter Richtárik
Cohort squeeze: Beyond a single communication round per cohort in cross-device federated learning
Oral at the NeurIPS 2024 Federated Learning Workshop
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

Prune at the clients, not the server: Accelerated sparse training in federated learnin 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

ΑΙ

The Next Generation Of Artificial Intelligence

Rob Toews Contributor ^① *I write about the big picture of artificial intelligence.*

Oct 12, 2020, 09:22pm EDT

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

https://www.forbes.com/sites/robtoews/2020/10/12/the-next-generation-of-artificial-intelligence/?sh=4d14f60159eb https://www.forbes.com/sites/robtoews/2020/10/29/the-next-generation-of-artificial-intelligence-part-2/?sh=e02f2567a304



Follow



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



The National Artificial Intelligence R&D Strategic Plan

Table of Contents

Ex	ecutive Summary	vii
Int	troduction to the National AI R&D Strategic Plan: 2023 Update	1
	AI as a National Priority	1
St	rategy 1: Make Long-Term Investments in Fundamental and Responsible AI Research	3
	Advancing Data-Focused Methodologies for Knowledge Discovery	3
	Fostering Federated ML Approaches	4
	Understanding Theoretical Capabilities and Limitations of AI	4
	Pursuing Research on Scalable General-Purpose AI Systems	5
	Developing AI Systems and Simulations Across Real and Virtual Environments	5
	Enhancing the Perceptual Capabilities of AI Systems	5
	Developing More Capable and Reliable Robots	6
	Advancing Hardware for Improved Al	6
	Creating AI for Improved Hardware	7
	Embracing Sustainable AI and Computing Systems	8

Federated Learning Issues & Tools



Part 2 Introduction



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



Three Types of Heterogeneity

Data	data distributions $\mathcal{D}_1, \ldots, \mathcal{D}_n$ can be different
Compute	compute times τ_1, \ldots, τ_n are nonzero and can be different
Communication	communication times $\theta_1, \ldots, \theta_n$ are nonzero and can be different

Typical Assumptions



Stochastic gradients have bounded variance:

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

Our Papers on Optimal Parallel SGD



Our Papers

First optimal parallel SGD under...

... computation (and/or data) heterogeneity

... communication (and computation) heterogeneity [Rennala SGD as a special case]

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

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

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





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

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*

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*

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*

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.



ELDEN RING



Optimal Parallel Stochastic Gradient Methods

	Data Heterogeneity $(\mathcal{D}_i \text{ different})$	Compute Heterogeneity $(\tau_i \text{ different})$	Communication Heterogeneity $(\theta_i \text{ different})$	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	×	\checkmark	 Image: A start of the start of	~		Inf	~	nearly
Amelie SGD Tyurin & R '24	~	 Image: A start of the start of	 Image: A start of the start of	✓		Inf	 Image: A start of the start of	~

Part 3 Previous Approaches to Parallelizing 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

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

Feng Niu leonn@cs.wisc.edu

Benjamin Recht Christopher Ré brecht@cs.wisc.edu 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 stateof-the-art performance on a variety of machine learning tasks. Several researchers have recently proposed schemes to parallelize SGD, but all require performancedestroying 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 for 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

NeurIPS 2020 Test of Time Award

Stephen Wright Cited by Department of Computer Sciences and Wisconsin Institute for Discovery, University of Wisconsin Verified email at cs.wisc.edu - Homepage Citations h-index Optimization i10-inde CITED BY Numerical Optimization (2nd edition 44606

J Nocadal, SJ Wright Springer Gradient projection for sparse reconstruction: Application to compressed sensing and other inverse problems MMT Figueirsch, RD weak, SJ Wright IEEE Journal of selected topics in signal processing 1 (4), 586-597	4365	2007	2017 2018 2019 2020 2021 2022 2023
Primal-dual interior-point methods SJ Wright Society for Industrial and Applied Mathematics	3629	1997	Public access
Hogwild: A lock-free approach to parallelizing stochastic gradient descent B Recht, C Ro, S Wright, F Niu Advance in Neural Information Bronzeline, Statemer, 503 701	2719	2011	0 articles
Sparse construction by separate approximation Standard Deviation Strends and S	2284	2009	not available Based on funding mandates

tions on signal processing 57 (7), 2479-2493

TITLE

Hogwil B Recht Advance

x	70 199	44 123 6000		
		4500	Hogwild: A lo	ock-free
8 2019 2020	0 2021 2022 2023	2024 0	Authors	Benjamin F
0 1010 101	5 1011 1011 1010	1014	Publication date	2011
access		VIEW ALL	Conference	Advances i
s lable		67 articles available	Pages	693-701

VIEW ALL

Since 2019

29504

approach to parallelizing stochastic gradient descent

Recht, Christopher Re, Stephen Wright, Feng Niu

in Neural Information Processing Systems

Description Stochastic Gradient Descent (SGD) is a popular algorithm that can achieve state-of-theart 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.

Total citations Cited by 2719



Scholar articles Hogwild!: A lock-free approach to parallelizing stochastic gradient descent B Recht, C Re, S Wright, F Niu - Advances in neural information processing systems, 2011

Cited by 2718 Related articles All 35 versions

Hogwild!: Alock□ freeapproach toparallelizingstochasticgradientdescent * RB NiuF - ... Systems. Granada, Spain, 2011 Cited by 2 Related articles

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-theart 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}) \coloneqq \mathbb{E}_{\xi \sim \mathcal{D}} f(\mathbf{x}; \xi) \},\$

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 ξ . 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 ∇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 $\xi_{k-1} \sim D$. Given M machines capable of computing these stochastic gradient estimates $\nabla [K; \xi]$ in parallel, one approach to parallelizing SGD is what we call "Minibath 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 $\mathbf{g}_{k-1}^{(1)}, \dots, \mathbf{g}_{k-1}^{(2)}$, and thet takes a minibatch SGD betweet, $\mathbf{x}_k = \mathbf{x}_{k-1} - \mathbf{v}_k \cdot \frac{1}{M} \sum_{m=1}^{M} \mathbf{g}_{m-1}^m$. This is a natural idea with long history [16, 18, 55] and it is a 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

36th Conference on Neural Information Processing Systems (NeurIPS 2022).

Sharper Convergence Guarantees for Asynchronous SGD for Distributed and Federated Learning

Anastasia Koloskova EPFL anastasia.koloskova@epfl.ch Sebastian U. Stich CISPA* EPFL 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 $\tau_{\rm max}$ and show that an ε -stationary point is reached after $O(\sigma^2 e^{-2} + \tau_{\rm max} e^{-1})$ iterations, where σ denotes the variance of stochastic gradients.

In this work we obtain (i) a tighter convergence rate of $O(\sigma^2 \varepsilon^{-2} + \sqrt{\tau_{max}\tau_{avg}} \varepsilon^{-1})$ without any change in the algorithm, where τ_{avg} is the average delay, which can be significantly smaller than τ_{max} . We also provide (ii) a simple delay-adaptive learning rate scheme, under which asynchronous SGD achieves a convergence rate of $O(\sigma^2 \varepsilon^{-2} + \tau_{avg} \varepsilon^{-1})$, and does not require any extra hyperparameter tuning nor extra communications. Our result allows to show for the first time that asynchronous SGD is advays fazer than mini-back. 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 [40,34], 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

36th Conference on Neural Information Processing Systems (NeurIPS 2022).

arXiv: June 16, 2022

arXiv: June 15, 2022

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

Peter Richtárik KAUST Saudi Arabia KAUST Saudi Arabia

Abstract

Parallelization is a popular strategy for improving the performance of iterative algorithms. Optimization methods are no exception: design of efficient parallel optimization methods and tight analysis of their throwcical properties are important research endeavors. While the minimax complexities are well known for sequential optimization methods, the theory of parallel optimization methods is ties sex splored and the sex set of the hade that have access to an unbiased stochastic gra ded variance. We consider a fixed computation model character each worker requiring a fixed but worker-dependent time to calculate stocha rove lower bounds and develop optimal algor

1 Introduction

We consider the none

 $\min_{x \in Q} \Big\{ f(x) := \mathbb{E}_{\xi \sim \mathcal{D}} \left[f(x; \xi) \right] \Big\},$ where $f : \mathbb{R}^d \times \mathbb{S}_{\xi} \to \mathbb{R}, Q \subseteq \mathbb{R}^d$, and ξ is a random variable with some distribution \mathcal{D} on \mathbb{S}_{ξ} . In machine learning, \mathbb{S}_{ξ} could be the space of all possible data, \mathcal{D} is the distribution of the training dataset, and $f(x, \xi)$ is the loss of a data sample ξ . In this paper we address the following natural setup:

(i) n workers are available to work in parallel, (ii) the ith worker requires τ_i seconds¹ to calculate a stochastic gradient of f.

The function f is L-smooth and lower-bounded (see Assumptions 7.1–7.2), and stochastic gradients are unbiased and σ^2 -variance-bounded (see Assumption 7.3).

PDF

1.1 Classical theory ¹Or any other unit of time.

In the nonconvex setting, gradient descent (GD) is an optimal method with respect to the number of gradient (∇f) calls (Lar., 2020; Nesterov, 2018; Carmon et al., 2020) for finding an approximately stationary point of f. Obviously, a key issue with GD is that it requires access to be exact gradient

37th Conference on Neural Information Processing Systems (NeurIPS 2023).

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 $(\mathcal{D}_i \text{ different})$	Compute Heterogeneity $(\tau_i \text{ different})$	Communication Heterogeneity $(\theta_i \text{ different})$	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	 Image: A start of the start of		Finite	×	big data regime
Freya SGD Tyurin, Gruntkowska & R '24	×	>	0	~		Finite	×	×
Fragile SGD Tyurin & R '24	×	~	<	~		Inf	 Image: A start of the start of	nearly
Amelie SGD Tyurin & R '24	~	~	~	~		Inf	 Image: A start of the start of	✓

Rennala SGD

Algorithmic idea: Minibatch SGD with asynchronous minibatch collection



Upper Bound



Matching Lower Bound

Theorem (informal)

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

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

seconds.

Rennala SGD = first optimal parallel SGD



Classical Oracle: Keeps Track of # Iterations



New Oracle: Keeps Track of Time



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



Figure 3: # of workers n = 10000.

The End (kind of)

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

Alexander Tyurin Peter Richtárik KAUST KAUST Saudi Arabia Saudi Arabia exandertiurin@gnail.com richtarik@gnail.

Abstract

Parallelization is a pepular strategy for impreving the performance of naturely approxima, Optimization methods are no conjective design of efficient parallel optimization methods and tight analysis of their theoretical properties are imposed to the performance of the performanc

1 Introduction

We consider the nonconvex optimization problem

$$\begin{split} & \underset{e \in \mathcal{R}}{\inf} \left\{ f(x) := \mathbb{E}_{\xi \sim \mathcal{D}} \left[f(x;\xi) \right] \right\}, \end{split} \tag{1}$$
 where $f: \mathbb{R}^d \times \mathbb{S}_{\xi} \to \mathbb{R}, Q \subseteq \mathbb{R}^d$, and ξ is a random variable with some distribution \mathcal{D} on \mathbb{S}_{ξ} . In machine learning, \mathbb{S}_{ξ} could be the space of all possible data. It is the distribution of the training dataset, and f(x) is the least of all distribution of the training dataset, and f(x) is the least of all distribution gamma terms.

(i) n workers are available to work in parallel,
 (ii) the ith worker requires τ_i seconds¹ to calculate a stochastic gradient of f.

The function f is L-smooth and lower-bounded (see Assumptions 7.1–7.2), and stochastic gradients are unbiased and σ^2 -variance-bounded (see Assumption 7.3).

PDF

1.1 Classical theory

In the nonconvex setting, gradient descent (GD) is an optimal method with respect to the number of gradient (V) f calls (Lan, 2020; Nesterov, 2018; Carmon et al., 2020) for finding an approximatel stationary point of I. Obviously, a key issue with GD is that it requires access to the exact gradient ¹Or any other unit of time.

37th Conference on Neural Information Processing Systems (NeurIPS 2023)

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 $(\mathcal{D}_i \text{ different})$	Compute Heterogeneity $(\tau_i \text{ different})$	Communication Heterogeneity $(\theta_i \text{ different})$	Smooth Nonconvex	Smooth Convex	Infinite / Finite Sum?	Supports Decentralized Setup?	Optimal Time Complexity?
Rennala SGD Tyurin & R (NeurIPS '23)	×	×	0	 Image: A start of the start of		Inf	×	~
Malenia SGD Tyurin & R (NeurIPS '23)	\bigcirc	~	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	×	V	0	~		Finite	×	big data regime
Freya SGD Tyurin, Gruntkowska & R '24	×	~	0	 Image: A second s		Finite	×	×
Fragile SGD Tyurin & R '24	×	×		~		Inf	 Image: A start of the start of	nearly
Amelie SGD Tyurin & R '24	~	~	v	×		Inf	~	✓

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

Malenia SGD



(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)	$ au_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 $(\mathcal{D}_i \text{ different})$	Compute Heterogeneity $(\tau_i \text{ different})$	Communication Heterogeneity $(\theta_i \text{ different})$	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		\bigcirc	Inf	×	~
Shadowheart SGD Tyurin, Pozzi, Ilin & R '24	×	~	 Image: A start of the start of	~		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	 Image: A start of the start of	nearly
Amelie SGD Tyurin & R '24	~	~	~	~		Inf	 Image: A start of the start of	~

Convex (Data Homogeneous) Regime

Method	Time Complexity
Minibatch SGD	$ au_n\left(\min\left\{rac{\sqrt{L}R}{\sqrt{arepsilon}},rac{M^2R^2}{arepsilon^2} ight\}+rac{\sigma^2R^2}{narepsilon^2} ight)$
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)	$\min_{m \in [n]} \left[\left(\frac{1}{m} \sum_{i=1}^{m} \frac{1}{\tau_i} \right)^{-1} \left(\min\left\{ \frac{\sqrt{LR}}{\sqrt{\varepsilon}}, \frac{M^2 R^2}{\varepsilon^2} \right\} + \frac{\sigma^2 R^2}{m\varepsilon^2} \right) \right]$
Lower Bound (Theorem B.4)	$\min_{m \in [n]} \left[\left(\frac{1}{m} \sum_{i=1}^{m} \frac{1}{\tau_i} \right)^{-1} \left(\min\left\{ \frac{\sqrt{LR}}{\sqrt{\varepsilon}}, \frac{M^2 R^2}{\varepsilon^2} \right\} + \frac{\sigma^2 R^2}{m\varepsilon^2} \right) \right]$
Lower Bound (Section M) (Woodworth et al., 2018)	$\tau_1 \min\left\{\frac{\sqrt{LR}}{\sqrt{\varepsilon}}, \frac{M^2 R^2}{\varepsilon^2}\right\} + \left(\frac{1}{n} \sum_{i=1}^n \frac{1}{\tau_i}\right)^{-1} \frac{\sigma^2 R^2}{n\varepsilon^2}$

 ∇f is *L*-Lipschitz, f is *M*-Lipschitz, and $||x^0 - x^*|| \leq R$

Further Extensions

Shadowheart SGD: Distributed Asynchronous SGD with Optimal T Complexity Under Arbitrary Computation and Communication Hetero

Nexander Tyurin⁺ Marta Pezzi⁺¹ Ivan Ilin⁺ Peter Richtárik⁺

Abstract We consider nonconvex stochastic optimization sublems in the occurrence controlling dis-	worker i to send a compressed message to the server; com- pression in performed via applying lossy communication compression to the communicated message (a vector from		
protectes in the approximate contracture au- soluted spray where the communication times from workers to a server can not be ignered, and the computation and communication times are potentially different for all workers. Using an unbiased compension technique, we develop a new method, Sadwahanter SDL, that results	R ⁴): see Def. 2.1; (d) the surver can broadcast compressed vectors to the workers in (at most) v _{mex} seconds, compression is per- formed via applying a lossy communication compression spenator to the communicated message (a vector from R ⁴) see Def. 8.1.		
improves the time complexities of all previous centralized methods. Mercerex, we show that the time complexity of Shadowheart SGD is op- timal in the family of centralized methods with compresend communication. We also consider the bidirectional songs, where broadcating from the source to the worken is non-negligible, and	The main goal of this work is to find an optimized repetitization strategy/method thut workd work uniformly well in all sor- naries characterized by the values of the computation times h_1, \dots, h_n and communication times r_1, \dots, r_n , and r_{max} , Since we allow these times is be arbitrarily heterogeneous designing a single algorithm that would be optimal in all these scenaries seems challenging.		
device a consponding method.	From the viewpoint of federated learning (Konchrj et al. 2006; Kainour et al., 2021), our work is a theoretical study of device heterogeneity. Moreover, our formalism captures both cours-sile and cross-device settings as special cases		
We consider the nonconvex smooth optimization problem $\min_{x \in \mathbb{R}^{d}} \left\{ f(x) := \mathbb{R}_{\mathbb{C} \sim T_{\mathbb{T}}} \left\{ f(x; \xi) \right\}, (1)$	Due to our in-depth focus on device heterogeneity and the challenges that need to be evercome, we do not consider statistical heterogeneity, and leave an extension to this setup to future work.		
where $f(:\cdot) : \mathbb{R}^d \times S_{\xi} \rightarrow \mathbb{R}$, and D_{ξ} is a distribution on $S_{\xi} \neq \emptyset$. Given $c > 0$, we seek to find a possibility random point \hat{x} such that $\mathbb{E}[\nabla f(\hat{x}) ^2] \leq c$. Such a point \hat{x} is called	We rely on assumptions which are standard in the litera- ture on stochastic gradient methods: smoothness, lower boundedtess and bounded variance.		
an e-stationary point. We focus on solving the problem in the following setup:	Assumption 1.1. f is differentiable and L -smooth, i.e., $\ \nabla f(x) - \nabla f(y)\ \le L \ x - y\ , \forall x, y \in \mathbb{R}^d$.		
(a) n workers/nodes are able to compute abcohastic gradi- ents $\nabla f(z; \xi)$ of f_i is parallel and asynchronously, and it takes (at most) h_i seconds for worker (to compute a single stochastic evaluation	Assumption 1.2. There exist $f^* \in \mathbb{R}$ such that $f(x) \geq f^*$ for all $x \in \mathbb{R}^d$. We define $\Delta := f(x^0) - f^*$, where $x^0 \in \mathbb{R}^d$ is a starting point of all algorithms we consider.		
(b) the workers are connected to a server which acts as a commenciation hub; (c) the workers can communicate with the server in par- alist and asynchronously; it takes (at most) r ₁ seconds for	Assumption 1.3. For all $x \in \mathbb{R}^d$, the stochastic gradients $\nabla f(x;\xi)$ are unbiased, and their variance is bounded by $\sigma^2 \ge 0$, i.e., $\mathbb{R}[\nabla f(x;\xi)] = \nabla f(x)$ and $\mathbb{E}_{\mathbb{C}}[\nabla f(x;\xi) - \nabla f(x) ^2] \le \sigma^2$.		
¹ Kng Abdilá Usiversity of Science and Technology, Thereal, Sandi Azabia ¹ Usiversity of Pasia, Italy. Correspondence to: Alexander Tyurin <alexandertizein@gmail.com>.</alexandertizein@gmail.com>	To simplify the exposition, in what follows (up to Sec. 7) we first focus on the regime in which the broadcast cost can be ignored. We describe a strategy for extending our algorithm to the more energed regime in Sec. 8.		





Shadowheart SGD

Optimal Parallel SGD under Compute Heterogeneity & Communication Heterogeneity



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

Shadowheart SGD: Setup



Optimal Parallel Stochastic Gradient Methods

	Data Heterogeneity $(\mathcal{D}_i \text{ different})$	Compute Heterogeneity $(\tau_i \text{ different})$	Communication Heterogeneity $(\theta_i \text{ different})$	Smooth Nonconvex	Smooth Convex	Infinite / Finite Sum?	Supports Decentralized Setup?	Optimal Time Complexity?
Rennala SGD Tyurin & R (NeurIPS '23)	×	~	0	 Image: A start of the start of		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	×	~	\checkmark	~		Inf	×	~
Freya PAGE Tyurin, Gruntkowska & 7 24	×	~	0			Finite	×	big data regime
Freya S Tyurin, Grunt ⁺ wska & R '24	×	>	0	×		Finite	×	×
ragile SGD Tyurin & R '24	×	\checkmark	 Image: A second s	~		Inf	 Image: A start of the start of	nearly
Amelie SGD Tyurin & R '24	~	~	✓	×		Inf	 Image: A start of the start of	✓

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

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

Shadowheart SGD



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 RandK 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	$\begin{array}{l} \max\{h_n,\dot{\tau}_n\} \to \infty, \\ \max\{h_i,\dot{\tau}_i\} < \infty \forall i < n \\ (\text{the last worker is slow}) \end{array}$	Time Complexities in Some Regimes $h_i = h, \dot{ au}_i = \dot{ au} \ orall i \in [n]$ (equal performance)	Numeri 1	cal Compa $\sigma^2/\varepsilon = 10^3$	rison ^(b) 10 ⁶
Minibatch SGD (see (3))	$\max_{i \in [n]} \max\{h_i, d\dot{\tau}_i\} \left(\frac{L\Delta}{\varepsilon} + \frac{\sigma^2 L\Delta}{n\varepsilon^2} \right)$	∞ (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)	$\max_{i \in [n]} \max\{h_i, \dot{\tau}_i\} \left(\left(\frac{d}{n} + 1\right) \frac{L\Delta}{\varepsilon} + \frac{d\sigma^2 L\Delta}{n\varepsilon^2} \right)$	∞ (non-robust)	$\geq \frac{dh\sigma^2}{n\varepsilon} \frac{L\Delta}{\varepsilon}$ (worse, e.g., when ε small)	×3	$\times 10^2$	$\times 10^4$
Rennala SGD (Tyurin & Richtárik, 2023c), Asynchronous SGD (e.g., (Mishchenko et al., 2022))	$\geq \min_{j \in [n]} \max\left\{h_{\bar{\pi}_j}, d\dot{\tau}_{\bar{\pi}_j}, \frac{\sigma^2}{\varepsilon} \left(\sum_{i=1}^j \frac{1}{h_{\bar{\pi}_i}}\right)^{-1}\right\} \frac{L\Delta}{\varepsilon}^{(\mathbf{a})}$	$<\infty$ (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)	$\times 10^2$	×10	$\times 1.5$
Shadowheart SGD (see (9) and Alg. 1) (Corollary 4.4)	$t^*(d-1,\sigma^2/arepsilon,[h_i,\dot au_i]_1^n)rac{L\Delta}{arepsilon}^{(extbf{c})}$	$<\infty$ (robust)	$\max\left\{h,\dot{\tau},\frac{d\dot{\tau}}{n},\sqrt{\frac{d\dot{\tau}h\sigma^2}{n\varepsilon}},\frac{h\sigma^2}{n\varepsilon}\right\}\frac{L\Delta}{\varepsilon}$	×1	×1	×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 τ_i when apply Theorem N.5 because these methods send d coordinates. $\bar{\pi}$ is a permutation that sorts $\max\{h_i, d\dot{\tau}_i\} : \max\{h_{\bar{\pi}_1}, d\dot{\tau}_{\bar{\pi}_1}\} \le \cdots \le \max\{h_{\bar{\pi}_n}, d\dot{\tau}_{\bar{\pi}_n}\}$

^(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.



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

Shadowheart SGD: Adding More Workers...



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

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

202000

HARM 10 HAR

In real-world distributed syste encounter device betterogeneit units. These can stem from GP

> expectations frome, while others represented edges or even if all to participate in the training all the total edges of the energy of volving finite-sum nonconvex optimization problems of the form $\min_{\substack{n \in \mathbb{N}^d}} \left\{ f(n) := \frac{1}{n} \sum_{i=1}^m f(n) \right\},$ form $f_i : \mathbb{R}^d \to \mathbb{R}$ can be viewed as the bios of a machine learning model x on the i^{th} examples form g_i the $i^{th} \to \mathbb{R}$ can be viewed as the bios of a machine learning model x on the i^{th} examples. Our gail to bid an e -naturatory point, i.e., a (groudby)

it is view mat h_i(ψ) (f(x)) ≤ j, we become the transportational interpretation being:

there are a voorkardelinestablexizes able to work in parallel,
each worker than access to stochastic gradients ∇f_j j ∈ [m],
worker i calculates ∇f_j(j) in less or equal to v, ∈ [0, ∞] seconds for all i ∈ [n], j ∈ [m].

print. Under review.



Freya PAGE

Optimal Parallel SGD for Large-Scale Finite-Sum Problems



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*

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 $(\mathcal{D}_i \text{ different})$	Compute Heterogeneity $(\tau_i \text{ different})$	Communication Heterogeneity $(\theta_i \text{ different})$	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	×	~	 Image: A start of the start of	~		Inf	×	~
Freya PAGE Tyurin, Gruntkowska & R '24	×	 Image: A start of the start of	0	~		Finite	×	big data regime
Freya SGD Tyurin, Gruntkowska & R '24	×	 Image: A second s	0	~		Finite	×	×
Fragile SGD Tyurin & R '24	×	 Image: A start of the start of	 Image: A start of the start of	~		Inf	~	nearly
Amelie SGD Tyurin & R '24	~	~	 	~		Inf	~	~

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

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

PAGE: Optimal Serial SGD for Finite-Sum Nonconvex Optimization

PAGE: A Simple and Optimal Probabilistic Gradient Estimator for Nonconvex Optimization

Zhize Li¹ Hongyan Bao¹ Xiangliang Zhang¹ Peter Richtárik

In this paper, we propose a novel stochastic gradient estimator-ProbAbilistic Gradient Esti

Abstract

(Jain & Kar. 2017). Driven by the applied success of deer neural networks (LeCun et al., 2015), and the critical place nonconvex optimization plays in training them, research in nonconvex optimization has been undergoing a renais-(PAGE)-for nonconvex optimization, PAGE is sance (Ghadimi & Lan, 2013; Ghadimi et al., 2016; Zhou easy to implement as it is designed via a small adet al., 2018; Fang et al., 2018; Li, 2019; Li & Richtárik, justment to vanilla SGD: in each iteration, PAGE uses the vanilla minibatch SGD update with probability p_t or reuses the previous gradient with a 1.1. The problem small adjustment, at a much lower computational cost, with probability $1 - p_t$. We give a simple Motivated by this development, we consider the general formula for the optimal choice of p_t . Moreover, optimization problem we prove the first tight lower bound $\Omega(n + \frac{\sqrt{n}}{2})$ $\min_{x \in \mathbb{R}^d} f(x),$ where $f: \mathbb{R}^d \to \mathbb{R}$ is a differentiable and possibly non-

finite-sum form

for nonconvex finite-sum problems, which also leads to a tight lower bound $\Omega(b + \frac{\sqrt{b}}{c^2})$ for nonconvex online problems, where $b := \min\{\frac{\sigma^2}{2}, n\}$. Then, we show that PAGE obtains the optimal convergence results $O(n + \frac{\sqrt{n}}{\epsilon^2})$ (finite-sum) and $O(b + \frac{\sqrt{b}}{\epsilon^2})$ (online) matching our lower bounds for both nonconvex finite-sum and online problems. Besides, we also show that for nonconvex functions satisfying the Polyak-Łojasiewicz (PL) condition PAGE can automatically switch to a faster linear convergence rate $O(\cdot \log \frac{1}{\epsilon})$. Finally, we conduct several deep learning experiments (e.g., LeNet, VGG, ResNet) on real datasets in PyTorch showing that PAGE not only converges much faster than SGD in training but also achieves the higher test accuracy, validating the optimal theoretical results and confirming the practical superiority of PAGE.

1. Introduction

Nonconvex optimization is ubiquitous across many domains of machine learning, including robust regression, low rank matrix recovery, sparse recovery and supervised learning

¹King Abdullah University of Science and Technology, Thuwal, Kingdom of Saudi Arabia. Correspondence to: Zhize Li <zhize.li@kaust.edu.sa>. Proceedings of the 38th International Conference on Machine

Learning, PMLR 139, 2021. Copyright 2021 by the author(s).

risk minimization problems in machine learning (Shalev-Shwartz & Ben-David, 2014). Moreover, if the number of data samples n is very large or even infinite, e.g., in the online/streaming case, then f(x) usually is modeled via the online form $f(x) := \mathbb{E}_{\zeta \sim \mathcal{D}}[F(x, \zeta)],$ which we also consider in this work. For notational convenience, we adopt the notation of the finite-sum form (2) in the descriptions and algorithms in the rest of this paper However, our results apply to the online form (3) as well by letting $f_i(x) := F(x, \zeta_i)$ and treating n as a very large value or even infinite. 1.2. Gradient complexity

convex function. We are interested in functions having the

 $f(x) := \frac{1}{n} \sum_{i=1}^{n} f_i(x),$

where the functions f, are also differentiable and possi

bly nonconvex. Form (2) captures the standard empirical

To measure the efficiency of algorithms for solving the nonconvex optimization problem (1), it is standard to bound the number of stochastic gradient computations needed to find a solution of suitable characteristics. In this paper we

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

(1)

(2)

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

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

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

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

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

(after butchering/redefining notation)

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 ∇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)	$ au_1 m rac{\delta^0}{arepsilon} - ig(au_n rac{m}{n} rac{\delta^0}{arepsilon}ig)$	Suboptimal
Hero PAGE (Soviet PAGE) [Li et al., 2021]	$ au_1 m + au_1 rac{\delta^0}{arepsilon} \sqrt{m} \left(au_n rac{m}{n} + au_n rac{\delta^0}{arepsilon} rac{\sqrt{m}}{n} ight)$	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]	$\tfrac{\delta^0}{\varepsilon} \left(\left(\sum_{i=1}^n \tfrac{1}{\tau_i} \right)^{-1} \left(\tfrac{\sigma^2}{\varepsilon} + n \right) \right)$	Limitations: σ^2 -bounded variance assumption, suboptimal when ε is small.
Rennala SGD [Tyurin and Richtárik, 2023]	$\tfrac{\delta^0}{\varepsilon}\min_{j\in[n]}\left(\left(\sum_{i=1}^j \tfrac{1}{\tau_i}\right)^{-1}\left(\tfrac{\sigma^2}{\varepsilon}+j\right)\right)$	Limitations: σ^2 -bounded variance assumption, suboptimal when ε 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., $\sqrt{m} \ge n$ (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 ε 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 $\sqrt{m} \ge n$.

^(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 Δ , 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 { τ_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 $k = 0, 1, \ldots, K - 1$ do 3: $x^{k+1} = x^k - \gamma q^k$ Sample $c^k \sim \text{Bernoulli}(p)$ 4: if $c^k = 1$ then 5: (with probability p) $\nabla f(x^{k+1}) = \text{ComputeGradient}(x^{k+1})$ 6: (Alg. 2) $g^{k+1} = \nabla f(x^{k+1})$ 7: 8: else (with probability 1-p) $\frac{1}{S}\sum_{i\in\mathcal{S}^k} \left(\nabla f_i(x^{k+1}) - \nabla f_i(x^k)\right) = \text{ComputeBatchDifference}(S, x^{k+1}, x^k) \quad \text{(Alg. 3)}$ 9: $g^{k+1} = g^k + \frac{1}{S} \sum_{i \in \mathcal{S}^k} \left(\nabla f_i(x^{k+1}) - \nabla f_i(x^k) \right)$ 10: end if 11: 12: **end for**

(note): S^k is a set of i.i.d. indices that are sampled from [m], uniformly with replacement, $|S^k| = S$

Algorithm 2 ComputeGradient(x)			Algorithm 3 ComputeBatchDifference (S, x, y)		
1:	Input: point $x \in \mathbb{R}^d$	1:	Input: batch size $S \in \mathbb{N}$, points $x, y \in \mathbb{R}^d$		
2:	Init $g = 0 \in \mathbb{R}^d$, set $\mathcal{M} = \emptyset$	2:	Init $g = 0 \in \mathbb{R}^d$		
3:	Broadcast x to all workers	3:	Broadcast x, y to all workers		
4:	For each worker $i \in [n]$, sample j from $[m]$	4:	For each worker, sample j from $[m]$ uniformly		
	uniformly and ask it to calculate $\nabla f_i(x)$		and ask it to calculate $\nabla f_i(x) - \nabla f_i(y)$		
5:	while $\mathcal{M} \neq [m]$ do	5:	for $i = 1, 2,, S$ do		
6:	Wait for $\nabla f_p(x)$ from a worker	6:	Wait for $\nabla f_p(x) - \nabla f_p(y)$ from a worker		
7:	if $p \in [m] ackslash \mathcal{M}$ then	7:	$g \leftarrow g + \frac{1}{S} (\nabla f_p(x) - \nabla f_p(y))$		
8:	$g \leftarrow g + rac{1}{m} abla f_p(x)$	8:	Sample j from $[m]$ uniformly and ask		
9:	Update $\mathcal{M} \leftarrow \mathcal{M} \cup \{p\}$		this worker to calculate $\nabla f_i(x) - \nabla f_i(y)$		
10:	end if	9:	end for		
11:	Sample j from $[m] \setminus \mathcal{M}$ uniformly and ask	10:	Return g		
	this worker to calculate $\nabla f_i(x)$				
12:	end while	Note	s: 1) the workers can aggregate ∇f_p locally, and the algorithm can		
12.	Poturn $a = \frac{1}{\sum_{n=1}^{m} \nabla f(x)}$	call	AllReduce once to collect all calculated gradients. ii) By splitting		
13:	Ketuin $g = \frac{1}{m} \sum_{i=1}^{N} \sqrt{J_i(x)}$	[m]	into blocks, instead of one $ abla f_p$, we can ask the workers to calculate		
	i=1	than	y_{m} of one block in Alg. 2 (and y_{m} a similar idea in Alg. 2)		

the sum of one block in Alg. 2 (and use a similar idea in Alg. 3).

Freya PAGE: Experiment 1



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

On the Optimal Time Complexities in Decentra Stochastic Asynchronous Ontimization

Alexander Tyurin Peter Bichtärik ing Abdallah University of Science and Technology (KALIST) Sandi Arabia (alexandertiurio, richtarik)dgaall.com

ic consider the decentralized stochastic asynchronous optimization setup, wh any workness asynchronouly calculate stochastic gradients and asynchronous municates with each other using degits in a miligraph, For both homogenees ab deterogeneous setup, we prove new time complexity lower bounds under asynchron the comparison and commissionic speeds an bounded. We devel new sampton that compared and homogeneous compatitiona and communication, that other and here the straight SGL and a new optimal method, here are sample optimal method. Pragits SGL and a new optimal method, here (50, the covery part of herbury here programs compatition and communication).

1 Introduction

 $\min_{x \in \mathbb{R}^d} \left\{ f(x) := \right.$

want to find a possibly random point z, called an c-stationary point, such that $\mathbb{E}[|\nabla f(z)|$ We analyze the hoterogeneous scrup and the corvex setup with smooth and non-smooth fund Sections B and C.

L1 Decentralized setup with th

we dividually use a second gravitation of the dividual weak, we can be a set of the second γ weak of the second γ and the second γ are associated as the second γ and the second γ and the second γ are associated as the second γ and the second γ and the second γ are associated as the second γ and the second γ are associated as the second γ and the second γ are associated as the second γ and the second γ are associated as the second γ and the second γ are associated as the second γ and the second γ are associated as the second γ are associated as the second γ and the second γ are associated as the second γ and the second γ are associated as the second γ are aspeciated as the second γ are associated as the secon

are dynamic. We consider any weighted directed multipuph parameterized by a vector $h \in \mathbb{R}^n$ such the $[0, \infty]$, and a matrix of distances $(p_{n+1})_{i,j} \in \mathbb{R}^{n \times n}$ such that $p_{n+j} \in [0, \infty]$ for all $i, j \in [$ $n \to 0$ for all $i \in [n]$. There were the intermediated are other surder.

Amelie SGD

Optimal Decentralized SGD under Computation & Communication Heterogeneity



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

Decentralized Setup: Amelie SGD

Method	The Worst-Case Time Complexity Guarantees	Comment
Minibatch SGD	$\frac{L\Delta}{\varepsilon} \max\left\{ \left(1 + \frac{\sigma^2}{n\varepsilon}\right) \max\{\max_{i,j\in[n]} \tau_{i\to j}, \max_{i\in[n]} h_i\} \right\}$	suboptimal if σ^2/ε 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 σ^2/ε 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)	$\frac{L\Delta}{\varepsilon} \max\left\{ \max_{i,j\in[n]} \tau_{i\to j}, \max_{i\in[n]} h_i, \frac{\sigma^2}{n\varepsilon} \left(\frac{1}{n} \sum_{i=1}^n h_i \right) \right\}$	Optimal up to a constant factor

The End (for real)