

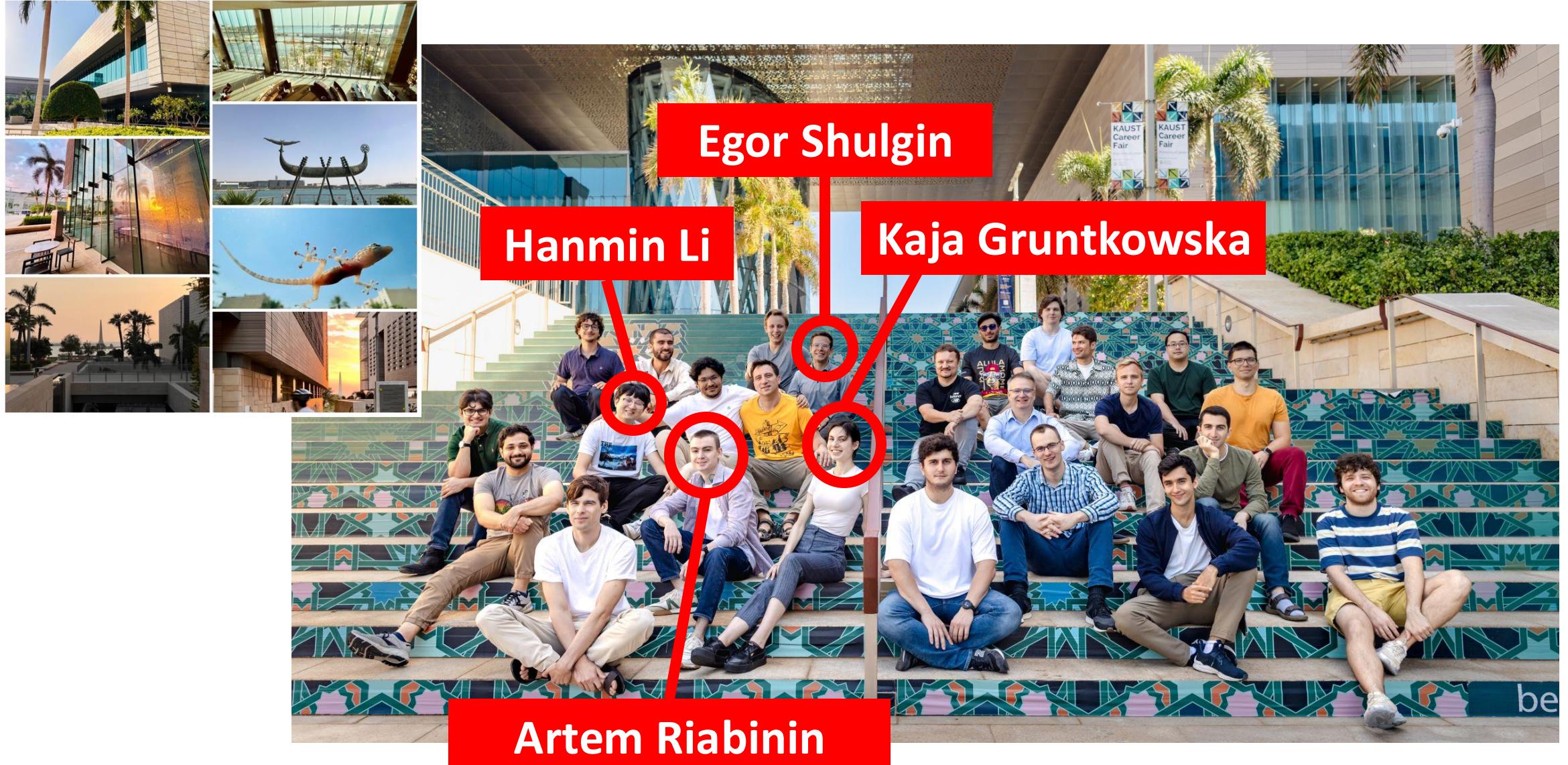
From the Broximal Point Method to Efficient Training of LLMs

Peter Richtárik
King Abdullah University of Science and Technology
Kingdom of Saudi Arabia



June 4-6, 2025 @ 3rd Vienna Workshop on Computational Optimization 2025 (VWCO25)

Optimization & Machine Learning Lab @ KAUST



Outline

Part 1: Broximal Point Method (BPM)

Part 2: Deep Learning Without Adam: Orthogonalizing the Gradient

Part 3: Gluon: Efficient Training of Large Language Models

The Ball-Proximal ("Broximal") Point Method:
a New Algorithm, Convergence Theory, and Applications

Kaja Gruntkowska¹, Hanmin Li¹, Aadi Rane^{*1,2}, Peter Richtárik¹

Abstract

1. Introduction

The minimization of non-convex global optimization poses significant challenges across various applications, where standard gradient-based methods often fail to find a global minimum. Such challenges have led to most modern learning algorithms, arising in both training and inference, where non-convex objectives or constraints are often necessary to capture complex prediction tasks.

1.1. Global nonconvex optimization

In this paper, we propose a new rate-of-convergence framework for global optimization problems, including non-convex and non-smooth optimization, based on smoothing, adaptive stepsize selection, and proximal operators. At the heart of our framework is the ball-proximal ("broximal") operator, which replaces the quadratic distance penalty by replacing the quadratic distance penalty

by a ball constraint. Surprisingly, and in sharp contrast to the well-known proximal point method (PPM) in the nonsmooth convex regime, we prove that BPM converges linearly and in a finite number of steps in the non-convex regime. Additionally, we show that the concept of ball-convexity, we prove that BPM retains the same global convergence guarantees under the same conditions as PPM, making it a useful tool for a broader class of potentially non-convex optimization problems. Just like PPM plays a key role of a conjugate method in the development and development of practically efficient algorithms and algorithms for non-convex optimization, we believe that the ball-convexity and the ball-proximal operator play a similar role in the development and development of practically efficient algorithms and algorithms for non-convex optimization.

1.2. The ball-proximal operator

A key inspiration for our method stems from the well-known Proximal Point Method (PPM) (Rockafellar, 1976), which iterates between a function $f : \mathbb{R}^d \rightarrow \mathbb{R} \cup \{-\infty\}$ (which is nonempty), closed, and lower semicontinuous. We let X_f be the set of all minimizers of f , and $f_* := \min_x f(x)$.

Definition 1.1 (Ball Proximal Operator). The ball proximal ("proximal") operator with radius $t > 0$ associated with a function $f : \mathbb{R}^d \rightarrow \mathbb{R} \cup \{-\infty\}$ is given by

$$\text{prox}_f^t(x) := \arg \min_{z \in \mathbb{R}^d} f(z), \quad (2)$$

where $B_t(x) := \{z \in \mathbb{R}^d : \|z - x\| \leq t\}$ and $\|\cdot\|$ is the standard Euclidean norm.

According to this definition, for a given input point $x \in \mathbb{R}^d$, $\text{prox}_f^t(x)$ is the set of all minimizers of f within the ball of radius t centered at x .

^{*}Work of Aadi Rane was conducted during a VFRP internship at KAUST.

¹KAUST

Part 1

Broximal Point Method



Kaja Gruntkowska, Hanmin Li, Aadi Rane, P.R.

The ball-proximal ("broximal") point method: a new algorithm, convergence theory, and applications
arXiv preprint arXiv:2502.02002, 2025

Problem Formulation

$$\min_{x \in \mathbb{R}^d} f(x)$$

f is proper, closed & convex

$$\text{Argmin } f \neq \emptyset$$

$$x_0 \in \text{dom } f$$

arXiv:2502.02002v1 [math.OC] 4 Feb 2025

The Ball-Proximal (“Broximal”) Point Method: a New Algorithm, Convergence Theory, and Applications

Kaja Gruntkowska¹ Hanmin Li¹ Aadi Rane^{* 1,2} Peter Richtárik¹

Abstract

Non-smooth and non-convex global optimization poses significant challenges across various applications, where standard gradient-based methods often struggle. We propose the *Ball-Proximal Point Method*, *Broximal Point Method*, or *Ball Point Method* (**BPM**) for short – a novel algorithmic framework inspired by the classical Proximal Point Method (**PPM**) (Rockafellar, 1976), which, as we show, sheds new light on several foundational optimization paradigms and phenomena, including non-convex and non-smooth optimization, acceleration, smoothing, adaptive stepsize selection, and trust-region methods. At the core of **BPM** lies the *ball-proximal (“broximal”)* operator, which arises from the classical proximal operator by replacing the quadratic distance penalty by a ball constraint. Surprisingly, and in sharp contrast with the sublinear rate of **PPM** in the nonsmooth convex regime, we prove that **BPM** converges *linearly* and in a *finite* number of steps in the same regime. Furthermore, by introducing the concept of ball-convexity, we prove that **BPM** retains the same global convergence guarantees under weaker assumptions, making it a powerful tool for a broader class of potentially non-convex optimization problems. Just like **PPM** plays the role of a conceptual method inspiring the development of practically efficient algorithms and algorithmic elements, e.g., gradient descent, adaptive step sizes, acceleration (Ahn & Sra, 2020), and “W” in AdamW (Zhuang et al., 2022), we believe that **BPM** should be understood in the same manner: as a blueprint and inspiration for further development.

1. Introduction

The minimization of non-convex functions is a fundamental challenge across many fields, including machine learning, optimization, applied mathematics, signal processing and operations research. Solving such problems is integral to most machine learning algorithms arising in both training and inference, where non-convex objectives or constraints are often necessary to capture complex prediction tasks.

1.1. Global nonconvex optimization

In this paper, we propose a new meta-algorithm (see Section 1.3) capable of finding *global* minimizers for a specific (new) class of non-convex functions. In particular, we introduce an algorithmic framework designed to solve the (potentially non-convex) optimization problem

$$\min_{x \in \mathbb{R}^d} f(x), \quad (1)$$

where $f : \mathbb{R}^d \mapsto \mathbb{R} \cup \{+\infty\}$ is assumed to be proper (which means that the set $\text{dom } f := \{x \in \mathbb{R}^d : f(x) < +\infty\}$ is nonempty), closed, and have at least one minimizer. We let \mathcal{X}_f be the set of all minimizers of f , and $f_* := \min_x f(x)$.

1.2. The ball-proximal operator

A key inspiration for our method stems from the well-known Proximal Point Method (**PPM**) (Rockafellar, 1976), which iteratively adds a quadratic penalty term to the objective (see Section 3 for more details) and solves a modified subproblem at each step. Building on this idea, we introduce the *ball-proximal (“broximal”)* operator:

Definition 1.1 (Ball-Proximal Operator). The *ball-proximal (“broximal”)* operator with radius $t > 0$ associated with a function $f : \mathbb{R}^d \mapsto \mathbb{R} \cup \{+\infty\}$ is given by

$$\text{brox}_f^t(x) := \arg \min_{z \in B_t(x)} f(z), \quad (2)$$

where $B_t(x) := \{z \in \mathbb{R}^d : \|z - x\| \leq t\}$ and $\|\cdot\|$ is the standard Euclidean norm.

According to this definition, for a given input point $x \in \mathbb{R}^d$, $\text{brox}_f^t(x)$ returns the minimizer(s) of f within the ball of radius t centered at x .

Ball Proximal (“Broximal”) Operator

From Proximal to Broximal Operator: From Penalty to Constraint

$$\text{prox}_f^\gamma(x) = \underset{y \in \mathbb{R}^d}{\operatorname{Arg\,min}} \left\{ f(y) + \frac{1}{2\gamma} \|y - x\|^2 \right\}$$

$$\text{brox}_f^t(x) = \underset{y \in \mathbb{R}^d}{\operatorname{Arg\,min}} \left\{ f(y) : \|y - x\| \leq t \right\}$$

Notation for the constraint:

$$\mathcal{B}(x, t) := \{y \in \mathbb{R}^d \mid \|y - x\| \leq t\}$$

penalty

constraint

The First Brox Theorem (GLRR 2025)

A1. $f : \mathbb{R}^d \rightarrow \mathbb{R} \cup \{+\infty\}$ is proper, closed and convex

A2. $x \in \text{dom } f$

$\text{brox}_f^t(x)$ is nonempty

$$\text{Argmin } f \cap \mathcal{B}(x, t) \neq \emptyset \quad \Rightarrow \quad \text{brox}_f^t(x) \subseteq \text{Argmin } f$$

$$\text{Argmin } f \cap \mathcal{B}(x, t) = \emptyset \quad \Rightarrow \quad \left\{ \begin{array}{l} \text{brox}_f^t(x) \text{ is a singleton} \\ \|\text{brox}_f^t(x) - x\| = t \end{array} \right.$$

The Second Brox Theorem (GLRR 2025)

A1. $f : \mathbb{R}^d \rightarrow \mathbb{R} \cup \{+\infty\}$ is proper, closed and convex

A2. $x \in \text{dom } f$

Let $u \in \text{brox}_f^{\textcolor{red}{t}}(x)$

There exists $c(x, \textcolor{red}{t}) \geq 0$ such that

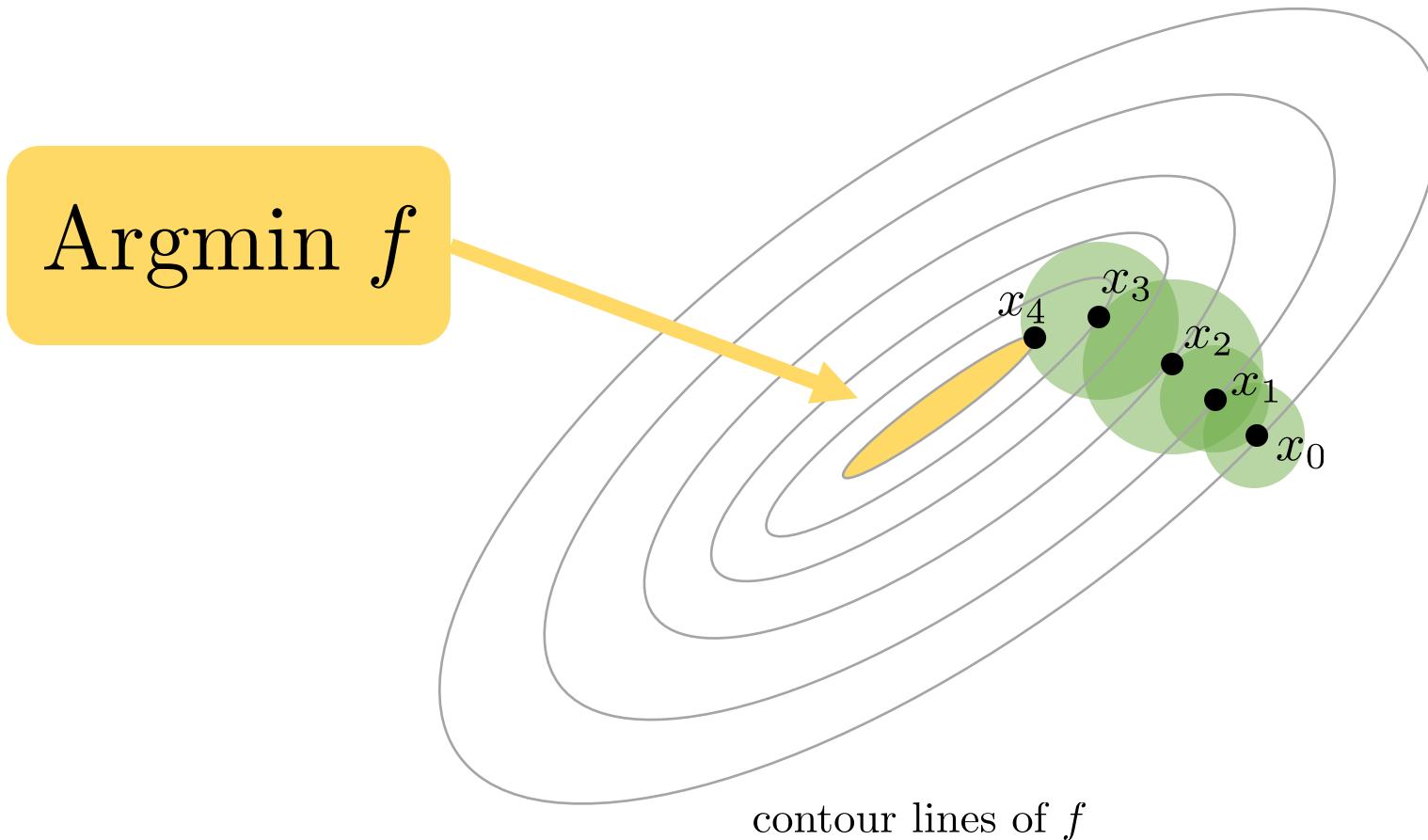
$$c(x, \textcolor{red}{t}) \cdot (x - u) \in \partial f(u)$$

$$f(y) - f(u) \geq c(x, \textcolor{red}{t}) \cdot \langle x - u, y - u \rangle \text{ for all } y \in \mathbb{R}^d$$

Broximal Point Method (BPM)

Broximal Point Method (BPM)

$$x_{k+1} = \text{brox}_f^{t_k}(x_k) = \operatorname{argmin} \{ f(x) : \|x - x_k\| \leq t_k \}$$



BPM = “Normalized” PPM

$f : \mathbb{R}^d \rightarrow \mathbb{R}$
convex & differentiable

$$x_{k+1} = \text{brox}_f^{t_k}(x_k) = x_k - t_k \frac{\nabla f(x_{k+1})}{\|\nabla f(x_{k+1})\|} = \text{prox}_f^{\gamma_k}(x_k)$$

$$\text{brox}_f^t(x) = \underset{y \in \mathbb{R}^d}{\text{Arg min}} \{f(y) : \|y - x\| \leq t\}$$

$$\gamma_k = \frac{t_k}{\|\nabla f(x_{k+1})\|}$$

normalized implicit gradient

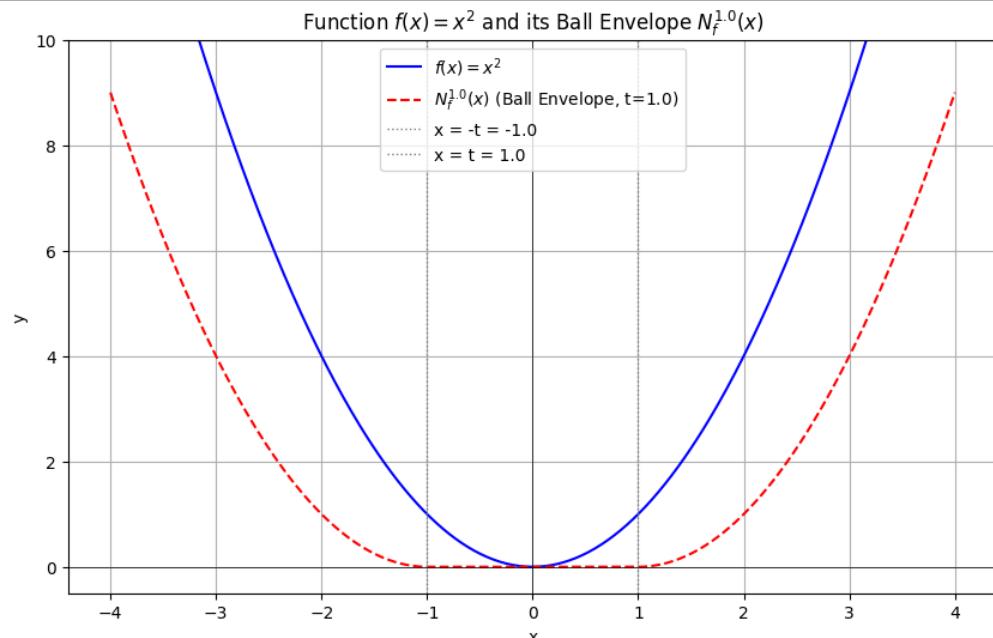
The Ball Envelope

Ball envelope of f : $N_f^t(x) := \min_{\|z-x\| \leq t} f(z)$

Argmin $N_f^t = \text{Argmin } f + \{x \in \mathbb{R}^d \mid \|x\| \leq t\}$

Example: $f(x) = x^2$

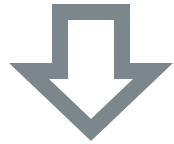
$$N_f^t(x) = \begin{cases} (x + t)^2 & \text{if } x < -t \\ 0 & \text{if } -t \leq x \leq t \\ (x - t)^2 & \text{if } x > t \end{cases}$$



BPM = Normalized GD for Minimizing the Ball Envelope

$f : \mathbb{R}^d \rightarrow \mathbb{R}$ is convex and L -smooth

$$x_{k+1} \notin \operatorname{Argmin} f$$



$$x_{k+1} = \text{brox}_f^{t_k}(x_k) = x_k - t_k \frac{\nabla N_f^{t_k}(x_k)}{\|\nabla N_f^{t_k}(x_k)\|}$$

Ball envelope of f :

$$N_f^t(x) := \min_{\|z-x\| \leq t} f(z)$$



Seven (Convergence) Results

Assumptions (Convex Setup)

f is proper, closed & convex

$\text{Argmin } f \neq \emptyset$

$x_0 \in \text{dom } f$

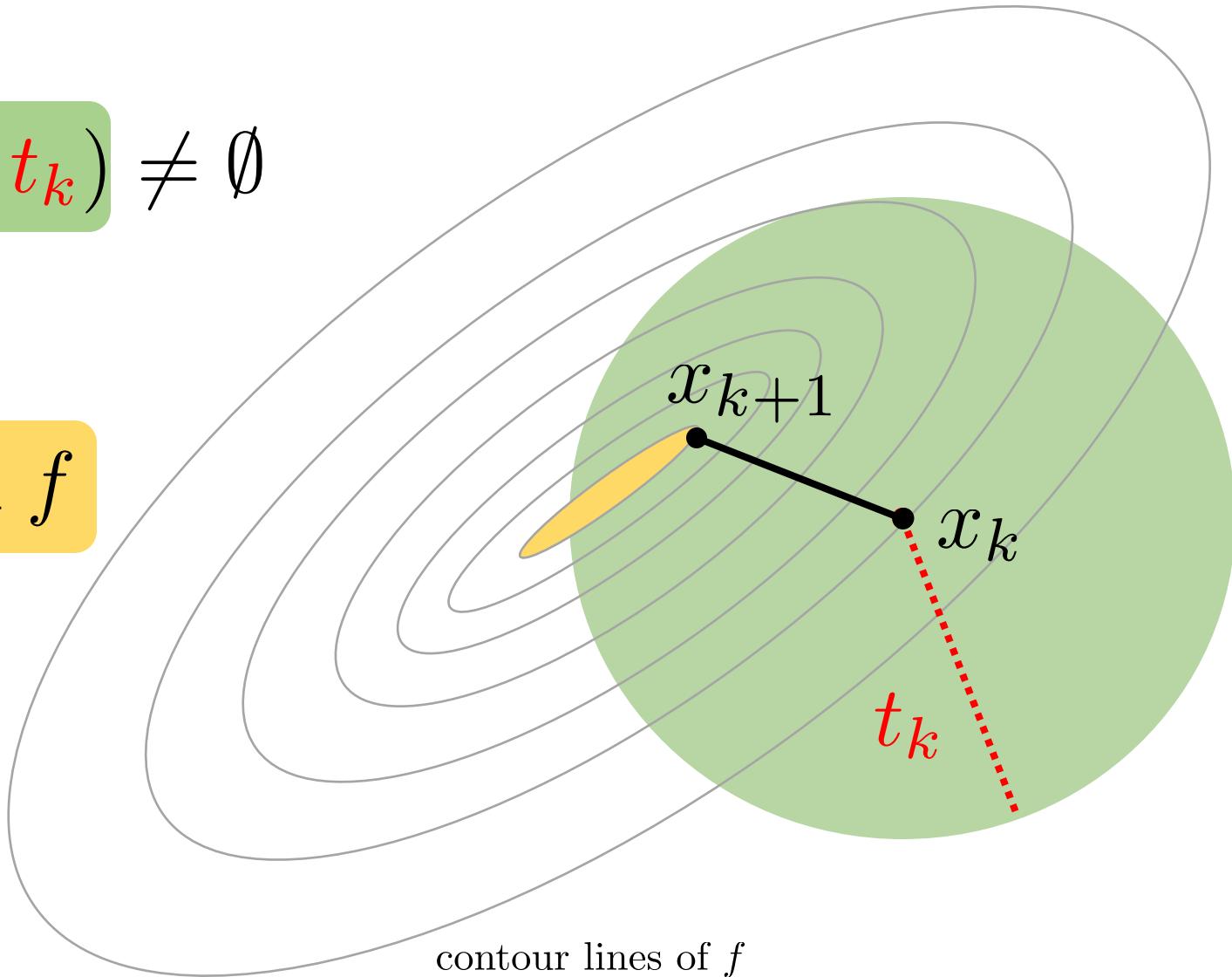
We have **nonconvex theory**
in the paper
(not covered in this talk)

Result 1: One Step Convergence

Argmin $f \cap \mathcal{B}(x_k, t_k) \neq \emptyset$



$x_{k+1} \in \text{Argmin } f$



Proximal Point Method (BPM)

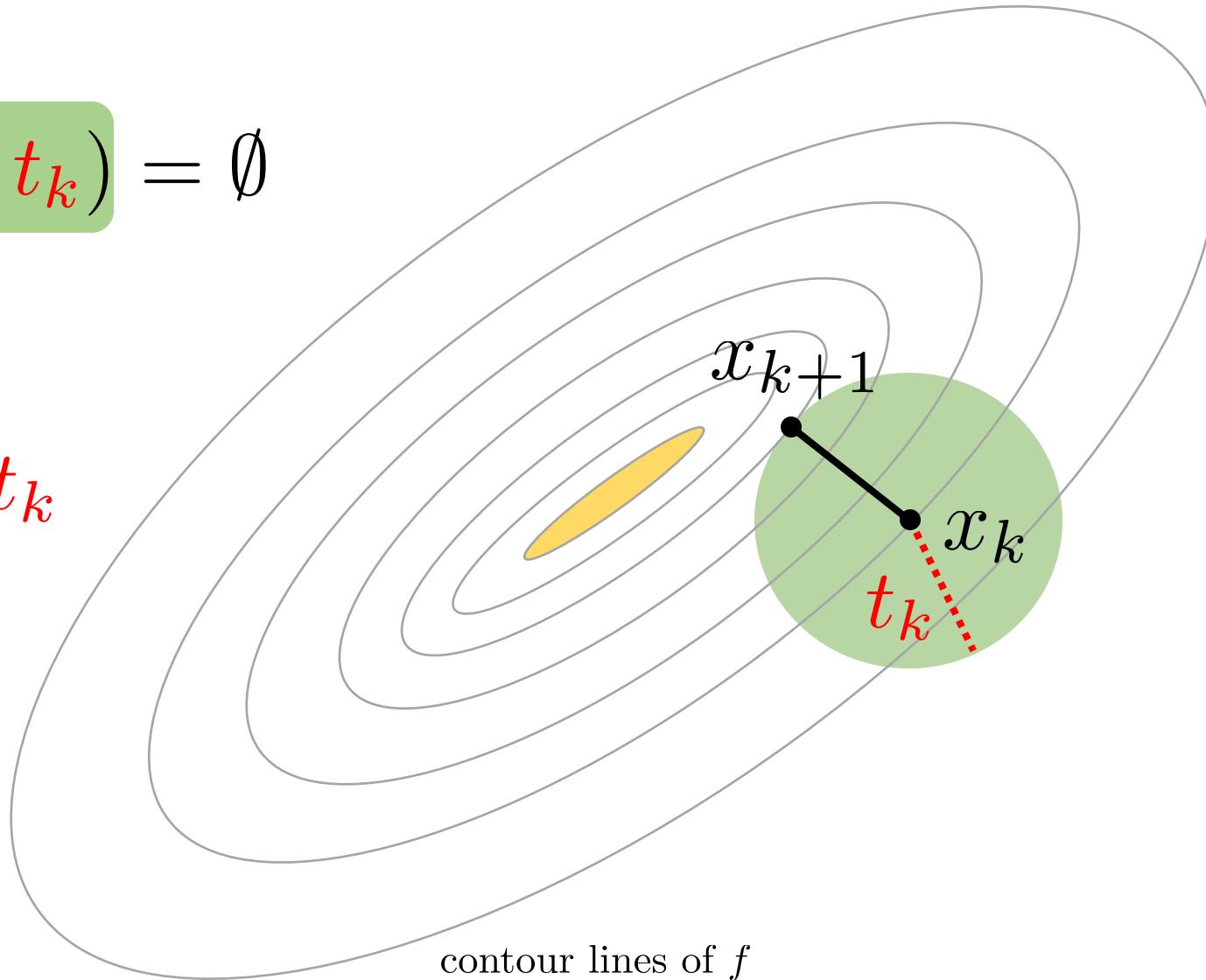
$$x_{k+1} \in \text{Argmin } \{f(x) : \|x - x_k\| \leq t_k\}$$

Result 2: Step to Boundary

$$\text{Argmin } f \cap \mathcal{B}(x_k, t_k) = \emptyset$$



$$\|x_{k+1} - x_k\| = t_k$$



Broximal Point Method (BPM)

$$x_{k+1} \in \text{Argmin } \{f(x) : \|x - x_k\| \leq t_k\}$$

Result 3: Super-Accelerated Linear Convergence in Distance to Minimizer

Argmin $f \cap \mathcal{B}(x_k, t_k) = \emptyset$

$x_\star \in \text{Argmin } f$



$$\|x_{k+1} - x_\star\|^2 \leq \|x_k - x_\star\|^2 - t_k^2$$

No need for
strong convexity!

Broximal Point Method (BPM)

$x_{k+1} \in \text{Argmin } \{f(x) : \|x - x_k\| \leq t_k\}$

$$\left(1 - \frac{t_k^2}{\|x_k - x_\star\|^2}\right) \|x_k - x_\star\|^2$$

Result 4: Finite Convergence (i.e., BPM is a “Direct” Optimization Method)

$$\sum_{k=0}^{K-1} t_k^2 \geq \|x_0 - x_\star\|^2 \quad \Rightarrow \quad x_K \in \text{Argmin } f$$

$x_\star \in \text{Argmin } f$

The diagram illustrates the finite convergence of the BPM. It starts with the statement $x_\star \in \text{Argmin } f$ in a yellow box at the top. A yellow arrow points down to the inequality $\sum_{k=0}^{K-1} t_k^2 \geq \|x_0 - x_\star\|^2$, which is followed by a large grey arrow pointing right to the conclusion $x_K \in \text{Argmin } f$ in a yellow box.

Corollary:

Constant stepsize regime: $t_k \equiv t$ for all $k \geq 0$

Broximal Point Method (BPM)

$$x_{k+1} \in \text{Argmin } \{f(x) : \|x - x_k\| \leq t\}$$

$$K \geq \frac{\|x_0 - x_\star\|^2}{t^2} \quad \Rightarrow \quad x_K \in \text{Argmin } f$$

Result 5: Super-Accelerated Linear Convergence in Function Values

follows from

$$f(x_{k+1}) \leq f(x_k) - t_k \frac{f(x_{k+1}) - f(x_\star)}{\|x_{k+1} - x_\star\|}$$

$$f(x_{k+1}) - f(x_\star) \leq \left(1 + \frac{t_k}{\|x_{k+1} - x_\star\|} \right)^{-1} (f(x_k) - f(x_\star))$$

Broximal Point Method (BPM)

$$x_{k+1} \in \operatorname{Argmin} \{f(x) : \|x - x_k\| \leq t_k\}$$



Result 6: Gradient Monotonicity

f is differentiable & $x_0 \in \text{dom } f$



$$\|\nabla f(x_{k+1})\| \leq \|\nabla f(x_k)\|$$

Broximal Point Method (BPM)

$$x_{k+1} \in \text{Argmin } \{f(x) : \|x - x_k\| \leq t_k\}$$

Result 7: Gradient Convergence

follows from

$$f(x_{k+1}) \leq f(x_k) - t_k \|\nabla f(x_{k+1})\|$$

$$x_\star \in \operatorname{Argmin} f$$

$$\sum_{k=0}^{K-1} \frac{t_k}{\sum_{s=0}^{K-1} t_s} \|\nabla f(x_{k+1})\| \leq \frac{f(x_0) - f(x_\star)}{\sum_{k=0}^{K-1} t_k}$$

Corollary:

Constant stepsize regime: $t_k \equiv t$ for all $k \geq 0$

$$\|\nabla f(x_{K+1})\| \leq \frac{1}{K} \sum_{k=0}^{K-1} \|\nabla f(x_{k+1})\| \leq \frac{f(x_0) - f(x_\star)}{t \cdot K}$$

Broximal Point Method (BPM)

$$x_{k+1} \in \operatorname{Argmin} \{f(x) : \|x - x_k\| \leq t_k\}$$

Practical Considerations

BPM: Practical Considerations

1. Some results extend to **arbitrary non-Euclidean norms** [GR 2025]
(choice of the norm is very important in deep learning, and training of LLMs!)
2. Can be made practical by **approximation**:

approximate f : apply BPM to a model of f

(gives rise to **non-Euclidean trust region methods**)

approximate brox_f^t : use an iterative technique, stop early

BPM: Related Works

Prior Related Work: Acceleration

arXiv:2003.08078v1 [math.OC] 18 Mar 2020

Acceleration with a Ball Optimization Oracle

Yair Carmon* Arun Jambulapati* Qijia Jiang* Yujia Jin* Yin Tat Lee[†]
Aaron Sidford* Kevin Tian*

Abstract

Consider an oracle which takes a point x and returns the minimizer of a convex function f in an ℓ_2 ball of radius r around x . It is straightforward to show that roughly $r^{-1} \log \frac{1}{\epsilon}$ calls to the oracle suffice to find an ϵ -approximate minimizer of f in an ℓ_2 unit ball. Perhaps surprisingly, this is not optimal: we design an accelerated algorithm which attains an ϵ -approximate minimizer with roughly $r^{-2/3} \log \frac{1}{\epsilon}$ oracle queries, and give a matching lower bound. Further, we implement ball optimization oracles for functions with locally stable Hessians using a variant of Newton's method. The resulting algorithm applies to a number of problems of practical and theoretical import, improving upon previous results for logistic and ℓ_∞ regression and achieving guarantees comparable to the state-of-the-art for ℓ_p regression.

1 Introduction

We study unconstrained minimization of a smooth convex objective $f : \mathbb{R}^d \rightarrow \mathbb{R}$, which we access through a *ball optimization oracle* $\mathcal{O}_{\text{ball}}$, that when queried at any point x , returns the minimizer¹ of f restricted a ball of radius r around x , i.e.,

$$\mathcal{O}_{\text{ball}}(x) = \arg \min_{x' \text{ s.t. } \|x'-x\| \leq r} f(x').$$

Such oracles underlie trust region methods [12] and, as we demonstrate via applications, encapsulate problems with local stability. Iterating $x_{k+1} \leftarrow \mathcal{O}_{\text{ball}}(x_k)$ minimizes f in $\tilde{O}(R/r)$ iterations (see Appendix A), where R is the initial distance to the minimizer, x^* , and $\tilde{O}(\cdot)$ hides polylogarithmic factors in problem parameters, including the desired accuracy.

Given the fundamental geometric nature of the ball optimization abstraction, the central question motivating our work is whether it is possible to improve upon this $\tilde{O}(R/r)$ query complexity. It is natural to conjecture that the answer is negative: we require R/r oracle calls to observe the entire line from x_0 to the optimum, and therefore finding a solution using less queries would require jumping into completely unobserved regions. Nevertheless, we prove that the optimal query complexity scales as $(R/r)^{2/3}$. This result has positive implications for the complexity for several key regression tasks, for which we can efficiently implement the ball optimization oracles.

1.1 Our contributions

Here we overview the main contributions of our paper: accelerating ball optimization oracles (with a matching lower bound), implementing them under Hessian stability, and applying the resulting techniques to regression problems.

^{*}Stanford University, {yairc,jambipati,qjiang2,yujiajin,sidford,kjtian}@stanford.edu.
[†]University of Washington, yintat@uv.edu.
¹In the introduction we discuss exact oracles for simplicity, but our results account for inexactness.

1

Prior Related Work: Non-Euclidean Balls

ℓ_∞ balls

Matrix Scaling and Balancing via Box Constrained Newton's Method and Interior Point Methods

Michael B. Cohen* Aleksander Madry† Dimitris Tsipras‡ Adrian Vladu‡
MIT MIT MIT MIT
micohen@mit.edu madry@mit.edu tsipras@mit.edu avladu@mit.edu

Abstract

In this paper, we study matrix scaling and balancing, which are fundamental problems in scientific computing, with a long line of work on them that dates back to the 1960s. We provide algorithms for both these problems that, ignoring logarithmic factors involving the dimension of the input matrix and the size of its entries, both run in time $\tilde{O}(m \log \kappa \log^2(1/\varepsilon))$ where ε is the amount of error we are willing to tolerate. Here, κ represents the ratio between the largest and the smallest entries of the optimal scalings. This implies that our algorithms run in nearly-linear time whenever κ is quasi-polynomial, which includes, in particular, the case of strictly positive matrices. We complement our results by providing a separate algorithm that uses an interior-point method and runs in time $\tilde{O}(m^{3/2} \log(1/\varepsilon))$.

In order to establish these results, we develop a new second-order optimization framework that enables us to treat both problems in a unified and principled manner. This framework identifies a certain generalization of linear system solving that we can use to efficiently minimize a broad class of functions, which we call *second-order robust*. We then show that in the context of the specific functions capturing matrix scaling and balancing, we can leverage and generalize the work on Laplacian system solving to make the algorithms obtained via this framework very efficient.

*This material is based upon work supported by the National Science Foundation under Grant No. 1111109 and Grant No. 1553428, and by the National Defense Science and Engineering Graduate Fellowship.
†This material is based upon work supported by the National Science Foundation under Grant No. 1553428.
‡This material is based upon work supported by the National Science Foundation under Grant No. 1111109 and Grant No. 1553428

arXiv:1704.02310v2 [cs.DS] 21 Aug 2017

ℓ_p balls

Convex optimization with p -norm oracles

Deeksha Adil Brian Bullins
Institute for Theoretical Studies Department of Computer Science
ETH Zürich Purdue University
deeksha.adil@eth-its.ethz.ch bbullins@purdue.edu

Arun Jambulapati Aaron Sidford
Department of Computer Science Department of Computer Science
University of Michigan Stanford University
jmlpati@umich.edu sidford@stanford.edu

November 1, 2024

Abstract

In recent years, there have been significant advances in efficiently solving ℓ_s -regression using linear system solvers and ℓ_2 -regression [Adil-Kyng-Peng-Sachdeva, J. ACM'24]. Would efficient ℓ_p -norm solvers lead to even faster rates for solving ℓ_s -regression when $2 \leq p < s$? In this paper, we give an affirmative answer to this question and show how to solve ℓ_s -regression using $\tilde{O}(n^{1+\nu})$ iterations of solving smoothed ℓ_s regression problems, where $\nu := \frac{1}{s} - \frac{1}{p}$. To obtain this result, we provide improved accelerated rates for convex optimization problems when given access to an $\ell_p^*(\lambda)$ -proximal oracle, which, for a point c , returns the solution of the regularized problem $\min_x f(x) + \lambda \|x - c\|_p^s$. Additionally, we show that the rates we establish for the $\ell_p^*(\lambda)$ -proximal oracle are near-optimal.

Contents

1 Introduction	2
2 General Proximal Point Algorithm	4
2.1 When $s < \infty$: Proximal Point Oracles	4
2.2 When $s = \infty$: Ball Oracle	8
3 Line-Search Free Method	10
4 Applications	14
4.1 ℓ_s -Regression	14
4.2 High-Order Smooth Optimization in General Norms	18
5 Lower Bounds	20
5.1 Proximal Point Oracles: $s < \infty$	20
5.2 When $s = \infty$: Ball Oracle	23

1

arXiv:2410.24158v1 [math.OC] 31 Oct 2024

Prior Related Work: Parallel Optimization

arXiv:2406.07373v1 [math.OC] 11 Jun 2024

Closing the Computational-Query Depth Gap in Parallel Stochastic Convex Optimization

Arun Jambulapati
University of Michigan
jmblpati@gmail.com

Aaron Sidford
Stanford University
sidford@stanford.edu

Kevin Tian
University of Texas at Austin
kjtian@cs.utexas.edu

Abstract

We develop a new parallel algorithm for minimizing Lipschitz, convex functions with a stochastic subgradient oracle. The total number of queries made and the query depth, i.e., the number of parallel rounds of queries, match the prior state-of-the-art, [CJJ⁺23], while improving upon the computational depth by a polynomial factor for sufficiently small accuracy. When combined with previous state-of-the-art methods, our result closes a gap between the best-known query depth and the best-known computational depth of parallel algorithms.

Our method starts with a *ball acceleration* framework of previous parallel methods, i.e., [CJJ⁺20, ACJ⁺21], which reduce the problem to minimizing a regularized Gaussian convolution of the function constrained to Euclidean balls. By developing and leveraging new stability properties of the Hessian of this induced function, we depart from prior parallel algorithms and reduce these ball-constrained optimization problems to stochastic unconstrained quadratic minimization problems. Although we are unable to prove concentration of the asymmetric matrices that we use to approximate this Hessian, we nevertheless develop an efficient parallel method for solving these quadratics. Interestingly, our algorithms can be improved using fast matrix multiplication and use nearly-linear work if the matrix multiplication exponent is 2.

Contents

1	Introduction	2
2	Technical overview	5
2.1	Framework: convolutions and acceleration	6
2.2	Hessian stability of the Gaussian convolution	8
2.3	Hessian optimization without Hessian approximation	10
2.4	Parallel maintenance of rank-one updates	10
3	Parallel optimization of quadratic subproblems	11
3.1	Composite stochastic optimization	12
3.2	Parallel maintenance of rank-1 updates	15
3.3	High-probability error bound reduction	19
3.4	Putting it all together	20
4	Parallel stochastic convex optimization	22
4.1	Approximate Newton's method	23
4.2	Ball optimization oracles via binary search	24
4.3	Proof of Theorem 2	28

ReSQueing Parallel and Private Stochastic Convex Optimization

Yair Carmon* Arun Jambulapati† Yujia Jin‡ Yin Tat Lee§ Daogao Liu†

Aaron Sidford‡ Kevin Tian§

Abstract

We introduce a new tool for stochastic convex optimization (SCO): a Reweighted Stochastic Query (ReSQue) estimator for the gradient of a function convolved with a (Gaussian) probability density. Combining ReSQue with recent advances in *ball oracle acceleration* [CJJ⁺20, ACJ⁺21], we develop algorithms achieving state-of-the-art complexities for SCO in parallel and private settings. For a SCO objective constrained to the unit ball in \mathbb{R}^d , we obtain the following results (up to polylogarithmic factors).

1. We give a parallel algorithm obtaining optimization error ϵ_{opt} with $d^{1/3}\epsilon_{\text{opt}}^{-2/3}$ gradient oracle query depth and $d^{1/3}\epsilon_{\text{opt}}^{-2/3} + \epsilon_{\text{opt}}^{-2}$ gradient queries in total, assuming access to a bounded-variance stochastic gradient estimator. For $\epsilon_{\text{opt}} \in [d^{-1}, d^{-1/4}]$, our algorithm matches the state-of-the-art oracle depth of [BJL⁺19] while maintaining the optimal total work of stochastic gradient descent.
2. Given n samples of Lipschitz loss functions, prior works [BFTT19, BFGT20, AFKT21, KLL21] established that if $n \gtrsim d\epsilon_{\text{dp}}^{-2}$, $(\epsilon_{\text{dp}}, \delta)$ -differential privacy is attained at no asymptotic cost to the SCO utility. However, these prior works all required a superlinear number of gradient queries. We close this gap for sufficiently large $n \gtrsim d^2\epsilon_{\text{dp}}^{-3}$, by using ReSQue to design an algorithm with near-linear gradient query complexity in this regime.

arXiv:2301.00457v2 [math.OC] 27 Oct 2023

*Tel Aviv University, ycarmon@tauex.tau.ac.il.

†University of Washington, {jmblpati, dgliu}@uw.edu.

‡Stanford University, {yujiajin, sidford}@stanford.edu.

§Microsoft Research, {yintatlee, tiankevin}@microsoft.com.

Prior Related Work: Minimize Maximum Loss

arXiv:2106.09481v2 [math.OC] 28 Oct 2021

Stochastic Bias-Reduced Gradient Methods

Hilal Asi Yair Carmon Arun Jambulapati Yujia Jin Aaron Sidford
 {asi,jmblpati,yujiajin,sidford}@stanford.edu ycarmon@cs.tau.ac.il

Abstract

We develop a new primitive for stochastic optimization: a low-bias, low-cost estimator of the minimizer x_* of any Lipschitz strongly-convex function. In particular, we use a multilevel Monte-Carlo approach due to Blanchet and Glynn [8] to turn any optimal stochastic gradient method into an estimator of x_* with bias δ , variance $O(\log(1/\delta))$, and an expected sampling cost of $O(\log(1/\delta))$ stochastic gradient evaluations. As an immediate consequence, we obtain cheap and nearly unbiased gradient estimators for the Moreau-Yoshida envelope of any Lipschitz convex function, allowing us to perform dimension-free randomized smoothing.

We demonstrate the potential of our estimator through four applications. First, we develop a method for minimizing the maximum of N functions, improving on recent results and matching a lower bound up to logarithmic factors. Second and third, we recover state-of-the-art rates for projection-efficient and gradient-efficient optimization using simple algorithms with a transparent analysis. Finally, we show that an improved version of our estimator would yield a nearly linear-time, optimal-utility, differentially-private non-smooth stochastic optimization method.

1 Introduction

Consider the fundamental problem of minimizing a μ -strongly convex function $F : \mathcal{X} \rightarrow \mathbb{R}$ given access to a stochastic (sub-)gradient estimator $\hat{\nabla}F$ satisfying $\mathbb{E}\hat{\nabla}F(x) \in \partial F(x)$ and $\mathbb{E}\|\hat{\nabla}F(x)\|^2 \leq G^2$ for every $x \in \mathcal{X}$. Is it possible to transform the unbiased estimator $\hat{\nabla}F$ into a (nearly) unbiased estimator of the minimizer $x_* := \operatorname{argmin}_{x \in \mathcal{X}} F(x)$? In particular, can we improve upon the $O(G/(\mu\sqrt{T}))$ bias achieved by T iterations of stochastic gradient descent (SGD)?

In this paper, we answer this question in the affirmative, proposing an *optimum estimator* \hat{x}_* , which (for any fixed $\delta > 0$) has

$$\text{bias } \|\mathbb{E}\hat{x}_* - x_*\| = \delta \text{ and variance } \mathbb{E}\|\hat{x}_* - \mathbb{E}\hat{x}_*\|^2 = O\left(\frac{G^2}{\mu^2} \log\left(\frac{G}{\mu\delta}\right)\right),$$

and, *in expectation*, costs $O(\log(\frac{G}{\mu\delta}))$ evaluations of $\hat{\nabla}F$.¹ Setting $\delta = G/(\mu\sqrt{T})$, we obtain the same bias bound as T iterations of SGD, but with expected cost of only $O(\log T)$ stochastic gradient evaluations (the worst-case cost is T). Further, the bias can be made arbitrarily small with only logarithmic increase in the variance and the stochastic gradient evaluations of our estimator, and therefore—paralleling the term “nearly linear-time” [27]—we call \hat{x}_* *nearly unbiased*.

Our estimator is an instance of the multilevel Monte Carlo technique for de-biasing estimator sequences [25] and more specifically the method of Blanchet and Glynn [8]. Our key observation is that this method is readily applicable to strongly-convex variants of SGD, or indeed any stochastic optimization method with the same (optimal) rate of convergence.

¹When $\mathcal{X} = \mathbb{B}_R(x_0) \subset \mathbb{R}^d$, $F(x) = \frac{1}{n} \sum_{i \in [n]} \hat{F}(x; i)$, and $\hat{\nabla}F$ is the subgradient of a uniformly random $\hat{F}(x; i)$ we can also get an estimator with bias 0 and expected cost $O(\log(nd))$. See Appendix A.1 for details.

arXiv:2105.01778v1 [math.OC] 4 May 2021

Thinking Inside the Ball:

Near-Optimal Minimization of the Maximal Loss

Yair Carmon Arun Jambulapati Yujia Jin Aaron Sidford
 [{ycarmon@cs.tau.ac.il}](mailto:ycarmon@cs.tau.ac.il), {jmblpati,yujiajin,sidford}@stanford.edu

Abstract

We characterize the complexity of minimizing $\max_{i \in [N]} f_i(x)$ for convex, Lipschitz functions f_1, \dots, f_N . For non-smooth functions, existing methods require $O(N\epsilon^{-2})$ queries to a first-order oracle to compute an ϵ -suboptimal point and $\tilde{O}(N\epsilon^{-1})$ queries if the f_i are $O(1/\epsilon)$ -smooth. We develop methods with improved complexity bounds of $O(N\epsilon^{-2/3} + \epsilon^{-8/3})$ in the non-smooth case and $O(N\epsilon^{-2/3} + \sqrt{N}\epsilon^{-1})$ in the $O(1/\epsilon)$ -smooth case. Our methods consist of a recently proposed ball optimization oracle acceleration algorithm (which we refine) and a careful implementation of said oracle for the softmax function. We also prove an oracle complexity lower bound scaling as $\Omega(N\epsilon^{-2/3})$, showing that our dependence on N is optimal up to polylogarithmic factors.

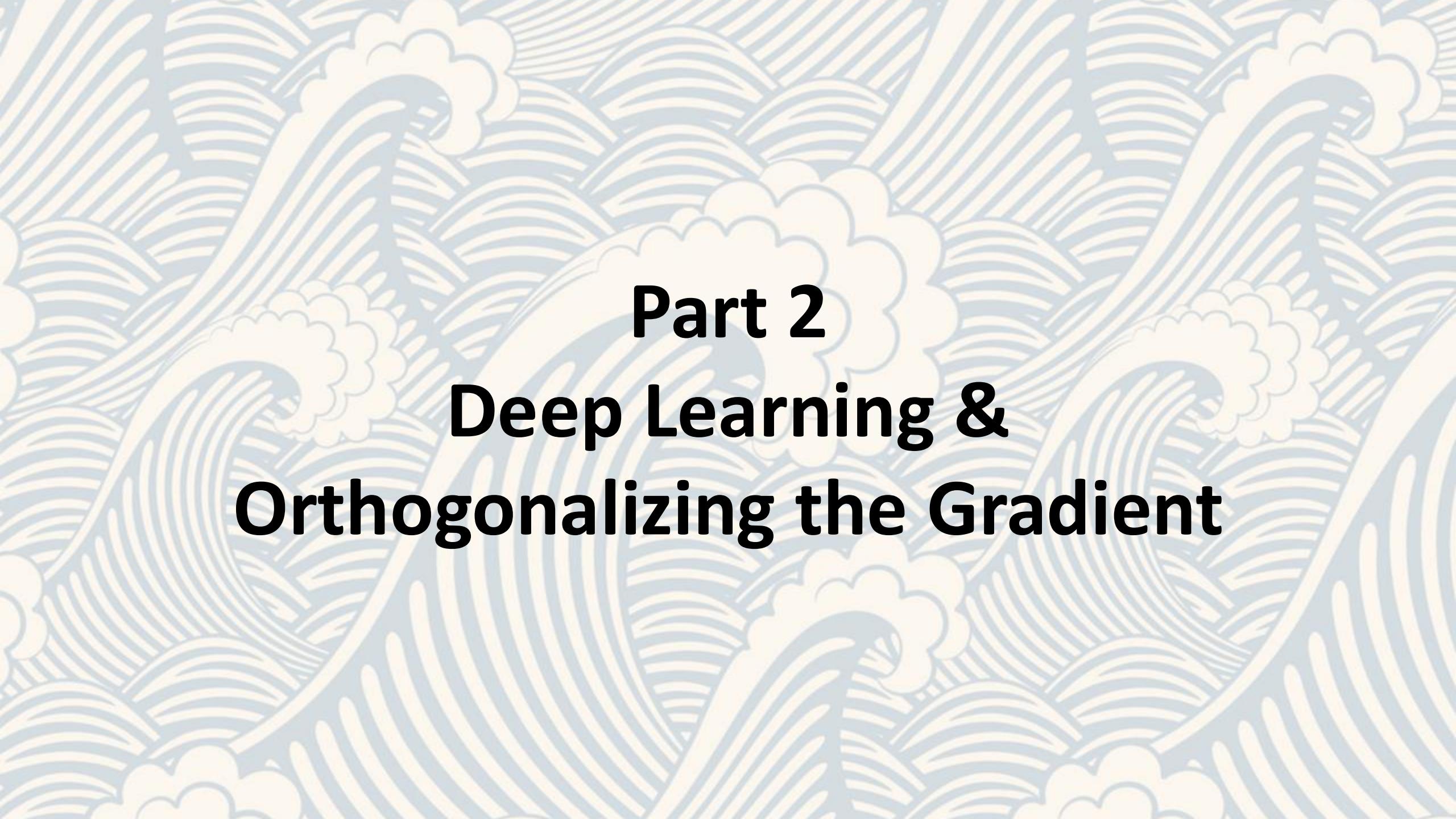
1 Introduction

Consider the problem of approximately minimizing the maximum of N convex functions: given f_1, \dots, f_N such that for every $i \in [N]$ the function $f_i : \mathbb{R}^d \rightarrow \mathbb{R}$ is convex, Lipschitz and possibly smooth, and a target accuracy ϵ ,

$$\text{find a point } x \text{ such that } F_{\max}(x) - \inf_{x_* \in \mathbb{R}^d} F_{\max}(x_*) \leq \epsilon \text{ where } F_{\max}(x) := \max_{i \in [N]} f_i(x). \quad (1)$$

Problems of this form play significant roles in optimization and machine learning. The maximum of N functions is a canonical example of structured non-smoothness and several works develop methods for exploiting it [31, 30, 36, 9, 12]. The special case where the f_i 's are linear functions is particularly important for machine learning, since it is equivalent to hard-margin SVM training (with f_i representing the negative margin on the i th example) [38, 13, 21]. Going beyond the linear case, Shalev-Shwartz and Wexler [36] argue that minimizing the maximum classification loss can have advantageous effects on training speed and generalization in the presence of rare informative examples. Moreover, minimizing the worst-case objective is the basic paradigm of robust optimization [4, 27]. In particular, since $F_{\max}(x) = \max_{p \in \Delta^N} \sum_{i \in [N]} p_i f_i(x)$ the problem corresponds to an extreme case of distributionally robust optimization [5] with an uncertainty set that encompasses the entire probability simplex Δ^N .

The goal of this paper is to characterize the complexity of this fundamental problem. We are particularly interested in the regime where the number of data points N and the problem dimension d are large compared to the desired level of accuracy $1/\epsilon$, as is common in modern machine learning. Consequently, we focus on dimension-independent first-order methods (i.e., methods which only rely on access to $f_i(x)$ and a (sub)gradient $\nabla f_i(x)$ as opposed to higher-order derivatives), and report complexity in terms of the number of function/gradients evaluations required to solve the problem.



Part 2

Deep Learning &

Orthogonalizing the Gradient

Shampoo (2018)

Algorithm 1 Shampoo, matrix case.

Initialize $W_1 = \mathbf{0}_{m \times n}$; $L_0 = \epsilon I_m$; $R_0 = \epsilon I_n$

for $t = 1, \dots, T$ **do:**

 Receive loss function $f_t : \mathbb{R}^{m \times n} \mapsto \mathbb{R}$

 Compute gradient $G_t = \nabla f_t(W_t)$ // $G_t \in \mathbb{R}^{m \times n}$

 Update preconditioners:

$$L_t = L_{t-1} + G_t G_t^\top$$

$$R_t = R_{t-1} + G_t^\top G_t$$

 Update parameters:

$$W_{t+1} = W_t - \eta L_t^{-1/4} G_t R_t^{-1/4}$$

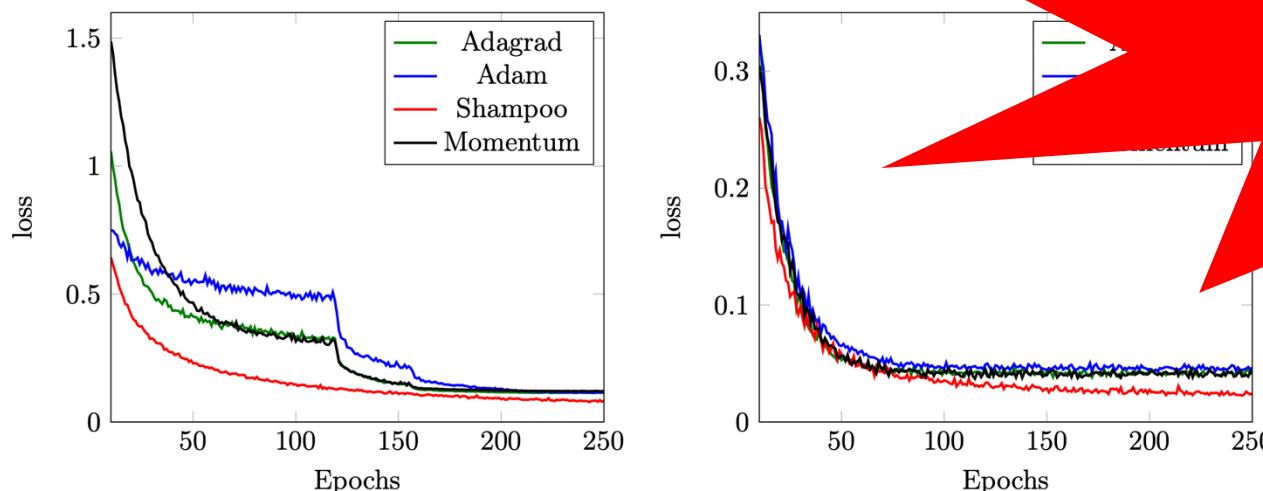


Figure 2: Training loss for a residual network and an inception network on CIFAR-10.

Shampoo: Preconditioned Stochastic Tensor Optimization

Vineet Gupta¹ Tomer Koren¹ Yoram Singer^{2,1}

Abstract

Preconditioned gradient methods are am-

most general and powerful tools in op-

timization.

However, preconditioning re-

quires manipulating prohibitively large

matrices. We propose to ob-

serve and analyze a new struc-

ture of the gradient in opti-

mization over tensors.

We propose a new algori-

thm, called Shampoo,

that uses a tensor precondi-

tioner to improve the con-

vergence properties of typi-

cal problems.

While preconditioned meth-

ods often lead to improved con-

vergence properties of typi-

cal problems.

However, the out-of-the-box use of full

preconditioners is ineffi-

cient and slow.

To mitigate this issue, specialized

algorithms have been de-

veloped in which the full precondi-

tioner is approximated by a

sketched version (Gonen

et al., 2014), a sketched version (Gonen

et al., 2017), or a

low-rank approximation (Nocedal,

1980) that can be used whenever second-order

information is unavailable or too expensive to compute.

Newer additions to this family are preconditioned online

algorithms, most notably AdaGrad (Duchi et al., 2011), that

use the covariance matrix of the accumulated gradients to

form a preconditioner.

While preconditioned methods often lead to improved con-

vergence properties of typi-

cal problems.

However, the out-of-the-box use of full

preconditioners is ineffi-

cient and slow.

To mitigate this issue, specialized

algorithms have been de-

veloped in which the full precondi-

tioner is approximated by a

sketched version (Gonen

et al., 2014), a sketched version (Gonen

et al., 2017), or a

low-rank approximation (Nocedal,

1980) that can be used whenever second-order

information is unavailable or too expensive to compute.

Newer additions to this family are preconditioned online

algorithms, most notably AdaGrad (Duchi et al., 2011), that

use the covariance matrix of the accumulated gradients to

form a preconditioner.

Promising, but did not
replace Adam
(259 citations only)

¹Google Brain, Mountain View, CA, USA ²Princeton University, Princeton, NJ, USA. Correspondence to: Tomer Koren <koren@google.com>.

Proceedings of the 35th International Conference on Machine Learning, Stockholm, Sweden, PMLR 80, 2018. Copyright 2018 by the author(s).

¹We call it Shampoo because it has to do with pre-conditioning.

Orthogonal SGDM (2022)

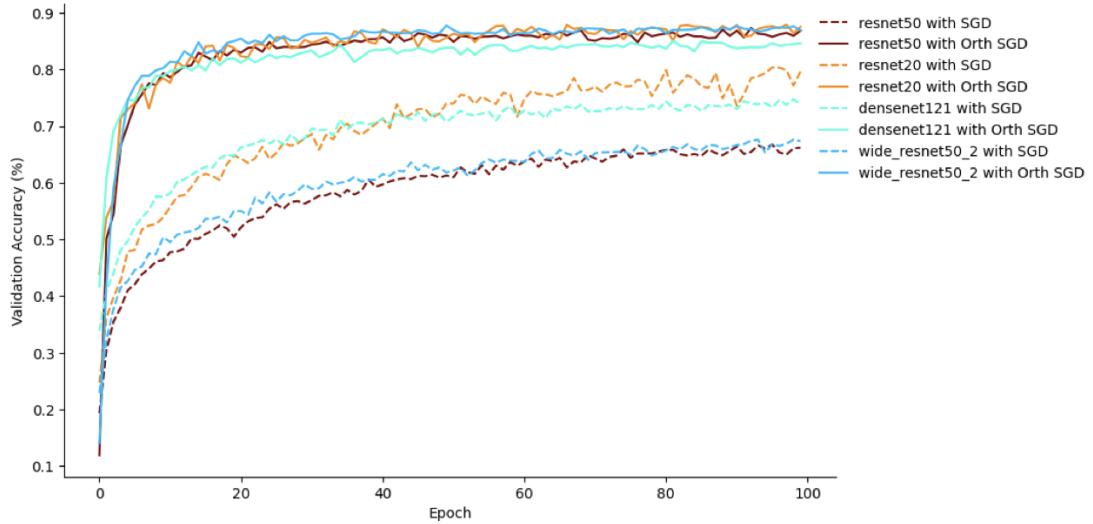


Figure 2: Validation accuracy from one run of SGDM vs Orthogonal-SGDM for a selection of models. Full plot in Appendix C. Best viewed in colour.

arXiv:2202.07052v1 [cs.LG] 14 Feb 2022

Orthogonalising gradients to speedup neural network optimisation

Mark Tuddenham¹ Adam Prügel-Bennett¹ Jonathon Hare¹

Abstract

The optimisation of neural networks can be sped up by orthogonalising the gradients before the optimisation step, ensuring the diversification of the learned representations. We orthogonalise the gradients of the layer's components/filters with respect to each other to separate out the intermediate representations. Our method of orthogonalisation allows the weights to be used more flexibly, in contrast to restricting the weights to an orthogonalised sub-space. We tested this method on ImageNet and CIFAR-10 resulting in a large decrease in learning time, and also obtain a speed-up on the semi-supervised learning BarlowTwins. We obtain similar accuracy to SGD without fine-tuning and better accuracy for naively chosen hyper-parameters.

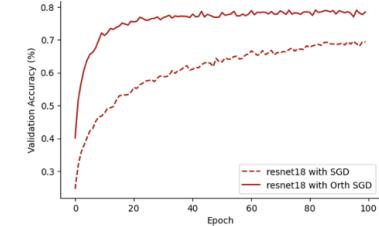


Figure 1: An example of the speed-up obtained by orthogonalising the gradients on CIFAR-10.

ties. Then, in Section 3, we provide experimental justifications and supporting experiments for this method along with finer details of the implementation and limitations.

2. Overview of new method and results

2.1. Related works

Gradient orthogonalisation has been explored in the domain of multi-task learning (Yu et al., 2020) to keep the different tasks separate and relevant. However in this work we focus on orthogonalisation for improving single task performance.

Weight orthogonalisation has been extensively explored with both empirical (Bansal et al., 2018; Jia et al., 2017) and theoretical (Jia et al., 2019) justifications. However, modifying the weights during training is unstable, and, in addition, it limits the weights to a tiny subspace. Deep learning is known to work well despite the immense size of the weight space, and as such we do not view this as an advantage. Xie et al. (2017) obtain improved performance over Stochastic Gradient Descent (SGD) via weight orthogonalisation and allows them to train very deep networks, we aim to achieve the same results while being more flexible with model and optimisation method choice. We do this by orthogonalising the gradients before they are used by an optimisation method rather than modifying the weights themselves. This, in effect, biases the training towards learning orthogonal representations and we show this to be the case;

¹Department of Electronics and Computer Science, University of Southampton, Southampton, UK. Correspondence to: Mark Tuddenham <mark.tuddenham@southampton.ac.uk>.

Preliminary work. Under review by the International Conference on Machine Learning (ICML). Do not distribute.

Orthogonal SGDM (2022)

2.2. Orthogonalising Gradients

Given a neural network f , with L layers made from components, c ,

$$f = \circ_{i=1}^L (f_i),$$
$$f_l(x) = \dots$$

where \circ is the composition of components in layer i , f_i is a parametrised function and $c_{l,i}$ denotes the parameters of f_i . Let $\theta_{l,i} \in \mathbb{R}^{P_l}$ giving $f_l : \mathbb{R}^{S_{l-1} \times N_{l-1}} \rightarrow \mathbb{R}^{P_l \times N_l}$ and let $\theta_l \in \mathbb{R}^{P_l \times N_l}$.

Let

$$G_l = [\nabla c_{l,1}, \nabla c_{l,2}, \dots, \nabla c_{l,N_l}], \quad (3)$$

be the $P_l \times N_l$ matrix of the components' gradients.

Then the nearest orthonormal matrix, *i.e.* the orthonormal matrix, O_l , that minimises the Frobenius norm of its difference from G_l

$$\min_{O_l} \|O_l - G_l\| \quad \text{subject to } \forall i, j : \langle O_{l,i}, O_{l,j} \rangle = \delta_{ij},$$

No theory
3 citations only as
of today!

where δ_{ij} is the Kronecker delta function, is the product of the left and right singular vector matrices from the Singular Value Decomposition (SVD) of G_l (Trefethen & Bau III, 1997),

$$G_l = U_l \Sigma_l V_l^\top, \quad (4)$$

$$O_l = U_l V_l^\top. \quad (5)$$

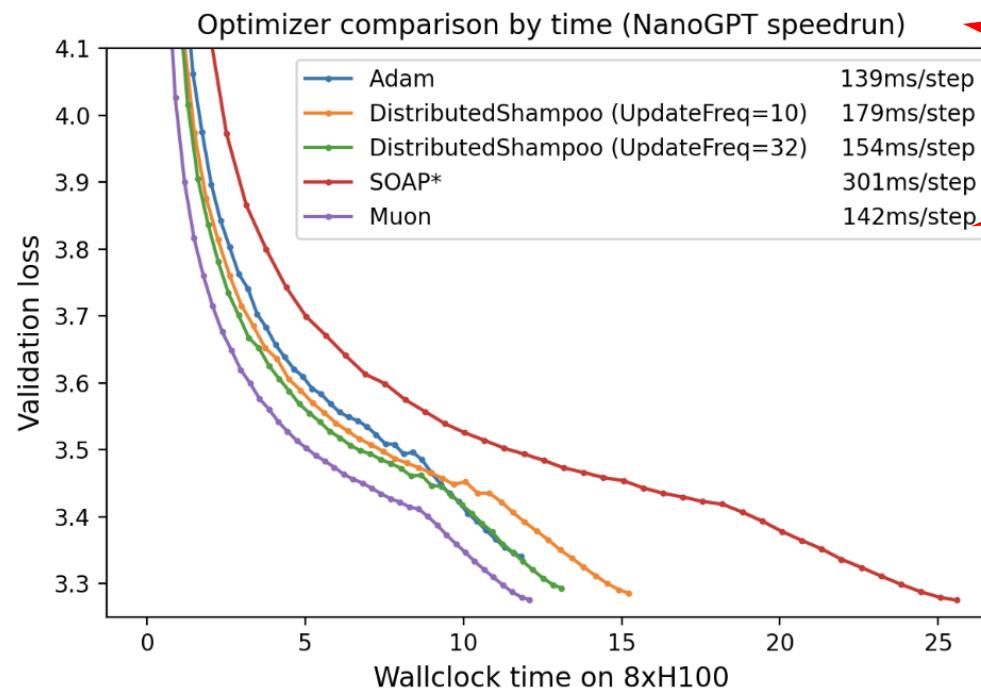
Adjust a first-order gradient descent method, Stochastic Gradient Descent with Momentum (Polyak, 1964), to make steps where the components are pushed in orthogonal directions,

$$v_l^{(t+1)} = \gamma v_l^{(t)} + \eta O_l^{(t)}, \text{ and} \quad (6)$$

$$\theta_l^{(t+1)} = \theta_l^{(t)} - v_l^{(t+1)}, \quad (7)$$

where v_l is the velocity matrix, $t \in \mathbb{Z}^{0+}$ is the time, γ is the momentum decay term, and η is the step size. We call this method Orthogonal Stochastic Gradient Descent with Momentum (Orthogonal-SGDM). This modification can clearly be applied to any first-order optimisation algorithm by replacing the gradients with $O_l^{(t)}$ before the calculation of the next iterate.

Muon (12/2024)



*SOAP is under active development. Future versions will significantly improve the wallclock overhead.

Figure 2. Optimizer comparison by wallclock time.

Muon: An optimizer for hidden layers in neural networks

December 8, 2024 · 16 min

Muon is an optimizer for the hidden layers of neural networks. It is used in the current training speedrun for both NanoGPT and AR-10 speedrunning.

Many experiments have been run to compare Muon to other optimizers. This writeup will provide an overview of Muon's performance and its potential for optimization.

Consistently beats Adam!

Algorithm 2 Muon

Require: Learning rate η , momentum μ

- 1: Initialize $B_0 \leftarrow 0$
- 2: **for** $t = 1, \dots$ **do**
- 3: Compute gradient $G_t \leftarrow \nabla_{\theta} \mathcal{L}_t(\theta_{t-1})$
- 4: $B_t \leftarrow \mu B_{t-1} + G_t$
- 5: $O_t \leftarrow \text{NewtonSchulz5}(B_t)$
- 6: Update parameters $\theta_t \leftarrow \theta_{t-1} - \eta O_t$
- 7: **end for**
- 8: **return** θ_t

Efficient Orthonormalization in Muon

Newton-Schulz 5

```
# Pytorch code
def newtonshulz5(G, steps=5, eps=1e-7):
    assert G.ndim == 2
    a, b, c = (3.4445, -4.7750, 2.0315)
    X = G.bfloat16()
    X /= (X.norm() + eps)
    if G.size(0) > G.size(1):
        X = X.T
    for _ in range(steps):
        A = X @ X.T
        B = b * A + c * A @ A
        X = a * X + B @ X
    if G.size(0) > G.size(1):
        X = X.T
    return X
```

By Keller Jordan (blog; 12/2024)

Optimal method (6/2025)

The Polar Express: Optimal Matrix Sign Methods and Their Application to the Muon Algorithm

Noah Amsel* David Persson† Christopher Musco‡ Robert M. Gower§

June 4, 2025

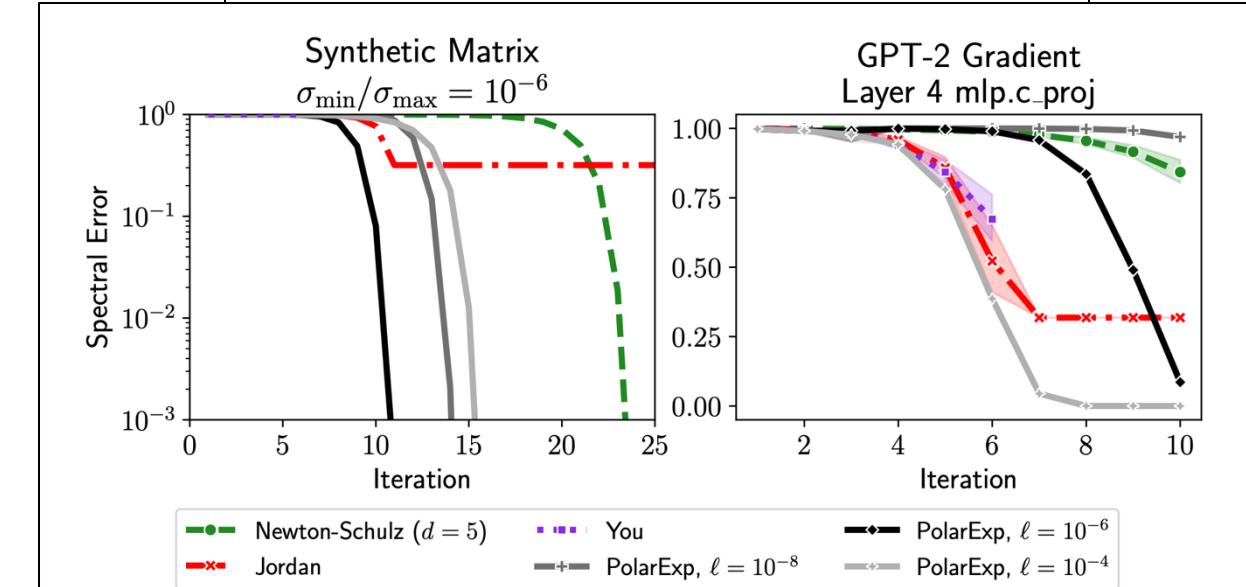


Figure 4: Convergence of various degree-5 polynomial methods in the spectral norm. When tuned properly, Polar Express attains outperforms the other methods at every iteration. Left panel: synthetic matrix with $\sigma_{\max} = 1$, $\sigma_{\min} = 10^{-6}$. Right panel: gradient of a certain weight matrix of a randomly-initialized GPT-2 architecture on a batch of language modeling data, normalized by the Frobenius norm.

Muon Theory

arXiv:2502.02900v2 [math.OC] 1 Jun 2025

A Note on the Convergence of Muon

Jiaxiang Li * Mingyi Hong †

Abstract

In this note, we inspect the convergence of a new optimizer for pretraining LLMs, namely the Muon optimizer. Such an optimizer is closely related to a specialized steepest descent method where the update direction is the minimizer of the quadratic approximation of the objective function under spectral norm (Bernstein and Newhouse, 2024). We provide the convergence analysis on both versions of the optimizer and discuss its implications.

1 Introduction

Recently, a new optimizer named Muon¹ has drawn attentions due to its success in training large models. Consider the following stochastic optimization problem:

$$\min_X f(X) = \mathbb{E}_\xi[F(x, \xi)] \quad (1.1)$$

where the variable $X \in \mathbb{R}^{m \times n}$ is a matrix. Without loss of generality we assume $m \geq n$. Muon optimizer writes:

$$\begin{aligned} G_t &\leftarrow \nabla F(X_t, \xi_t) \\ B_t &\leftarrow \mu B_{t-1} + G_t \\ O_t &\leftarrow \underset{\Delta}{\operatorname{argmin}} \{ \|O - B_t\|_F : O^\top O = I \text{ or } O O^\top = I\} \\ X_{t+1} &\leftarrow X_t - \eta_t O_t \end{aligned} \quad (1.2)$$

where the O step can be equivalently written as $O_t = U V^\top$ where $B_t = U \Sigma V^\top$ is the singular value decomposition.

On the other hand, a closely related formula of (1.2) is the following spectral optimizer:

$$\begin{aligned} G_t &\leftarrow \nabla F(X_t, \xi_t) \\ B_t &\leftarrow \mu B_{t-1} + G_t \\ \Delta_t &\leftarrow \underset{\Delta}{\operatorname{argmin}} \{ \operatorname{tr}(B_t^\top \Delta) + \frac{1}{2\eta_t} \|\Delta\|_2^2 \} \\ X_{t+1} &\leftarrow X_t + \Delta_t \end{aligned} \quad (1.3)$$

When $\|\cdot\|_2$ is the matrix 2-norm (spectral norm), i.e. the largest singular value, the solution of the third line of (1.3) is $\Delta_t = -\eta_t \|B_t\|_* U V^\top$, where $B_t = U \Sigma V^\top$ is the singular value decomposition (see Bernstein and Newhouse (2024, Proposition 5)). In this case (1.3) only differs from (1.2) by the norm $\|B_t\|_*$.

*Department of Electrical and Computer Engineering, University of Minnesota. li003755@umn.edu

†Department of Electrical and Computer Engineering, University of Minnesota. mhong@umn.edu

¹see <https://kellerjordan.github.io/posts/muon/>.

February 2025

UNDERSTANDING GRADIENT ORTHOGONALIZATION FOR DEEP LEARNING VIA NON-EUCLIDEAN TRUST-REGION OPTIMIZATION

A PREPRINT

Dmitry Kovalev
Yandex Research
dakovalev1@gmail.com

April 9, 2025

ABSTRACT

Optimization with matrix gradient orthogonalization has recently demonstrated impressive results in the training of deep neural networks (Jordan et al., 2024; Liu et al., 2025). In this paper, we provide a theoretical analysis of this approach. In particular, we show that the orthogonalized gradient method can be seen as a first-order trust-region optimization method, where the trust-region is defined in terms of the matrix spectral norm. Motivated by this observation, we develop the stochastic non-Euclidean trust-region gradient method with momentum, which recovers the Muon optimizer (Jordan et al., 2024) as a special case, along with normalized SGD and signSGD with momentum (Cukkosky and Mehta, 2020; Sun et al., 2023). In addition, we prove state-of-the-art convergence results for the proposed algorithm in a range of scenarios, which involve arbitrary non-Euclidean norms, constrained and composite problems, and non-convex, star-convex, first- and second-order smooth functions. Finally, our theoretical findings provide an explanation for several practical observations, including the practical superiority of Muon compared to the Orthogonal-SGD algorithm of Tuddenham et al. (2022) and the importance of weight decay in the training of large-scale language models.

1 Introduction

Over the past couple of decades, a substantial amount of optimization research has been dedicated to adaptive gradient optimization algorithms (Duchi et al., 2011; Tieleman, 2012; Kingma and Ba, 2014; Gupta et al., 2018; Reddi et al., 2019), with the primary application being the training of deep neural networks (LeCun et al., 2015). One of the most notable results of this line of research is the development of the AdamW optimizer (Loshchilov and Hutter, 2017; Zhuban et al., 2022), which has become the standard algorithm of choice for training large language models (LLMs) (Achiam et al., 2023; Liu et al., 2024; Grattafiori et al., 2024; Anil et al., 2023). However, recently, Jordan et al. (2024); Liu et al. (2025) have made significant progress in the ambitious task of surpassing AdamW in training LLMs using the idea of neural network optimization with orthogonalized gradients (Tuddenham et al., 2022; Jordan et al., 2024). In our paper, we aim to establish theoretical foundations for this promising research direction. Formally speaking, we consider the following composite optimization problem:

$$\min_{x \in \mathcal{X}} [f(x) = f(x) + R(x)], \quad (1)$$

where \mathcal{X} is a finite-dimensional vector space endowed with the inner product $\langle \cdot, \cdot \rangle : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$, $f(\cdot) : \mathcal{X} \rightarrow \mathbb{R}$ is a bounded from below and differentiable objective function, and $R(\cdot) : \mathcal{X} \rightarrow \mathbb{R} \cup \{-\infty\}$ is a proper,¹ closed, and convex regularizer. Our goal is to study the convergence properties of gradient methods for solving problem (1) in the *stochastic non-Euclidean smooth* setting. We provide a formal description of this setting below and justify our interest in this setting in the upcoming Sections 1.1 and 1.2.

Stochastic gradient estimator. We assume access to a stochastic estimator $g(\cdot; \xi) : \mathcal{X} \rightarrow \mathcal{X}$ of the gradient $\nabla f(\cdot)$, where $\xi \sim \mathcal{D}$ is a random variable sampled from a probability distribution \mathcal{D} . We assume that the stochastic gradient

¹Function $R(x)$ is called proper if there exists $x \in \mathcal{X}$ such that $R(x)$ is finite.

March 2025

Scion (2/2025)

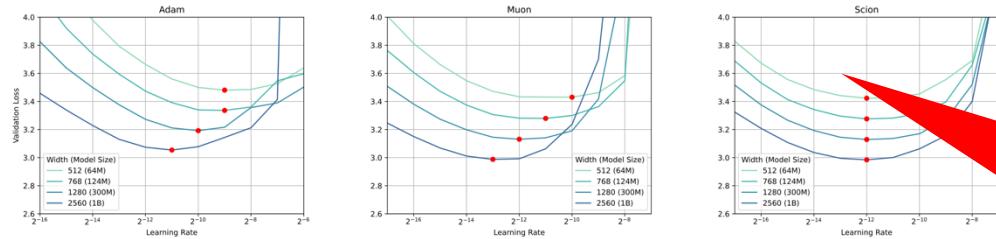


Figure 1. Performance on NanoGPT with between 64M and 1B parameters. The optimal learning rate of SCION is

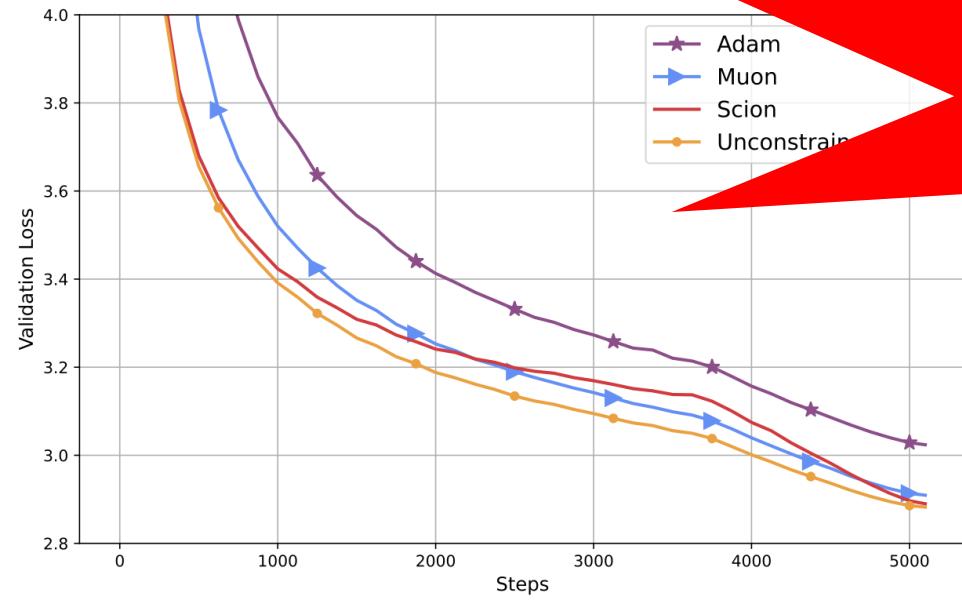


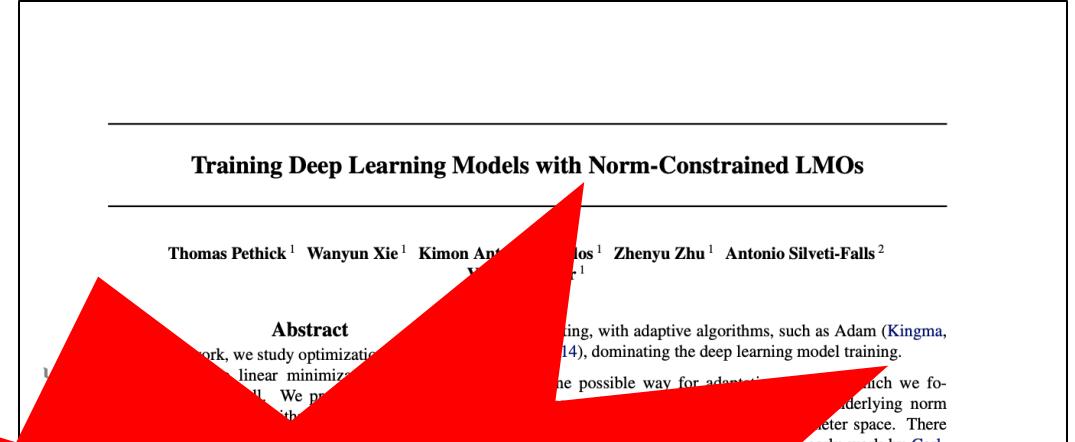
Figure 5. Validation loss curve for NanoGPT 3B.

Improves upon Muon,
exhibits hyperparameter
transfer

Input: $x^1, \dots, x^n \in \mathcal{D}$, $\alpha_k \in (0, 1)$, $\gamma_k \in (0, 1)$, $\xi_k \sim \mathcal{P}$, $d^0 \in \mathcal{D}$, $\text{lmo} : \mathcal{D} \rightarrow \mathcal{D}$
Output: $\bar{x}^n \in \mathcal{D}$

- 1: **for** $k = 1, \dots, n$ **do**
- 2: Sample $\xi_k \sim \mathcal{P}$
- 3: $d^k \leftarrow \alpha_k \nabla f(x^k, \xi_k) + (1 - \alpha_k)d^{k-1}$
- 4: $x^{k+1} \leftarrow (1 - \gamma_k)x^k + \gamma_k \text{lmo}(d^k)$
- 5: Choose \bar{x}^n uniformly at random from $\{x^1, \dots, x^n\}$

Return \bar{x}^n



Training Deep Learning Models with Norm-Constrained LMOs

Thomas Pethick¹ Wanyun Xie¹ Kimon Antonakos¹ Vassilios Kotsikas¹ Zhenyu Zhu¹ Antonio Silveti-Falls²

Abstract

In this work, we study optimization for deep learning models with adaptive algorithms, such as Adam (Kingma, 2014), dominating the deep learning model training.

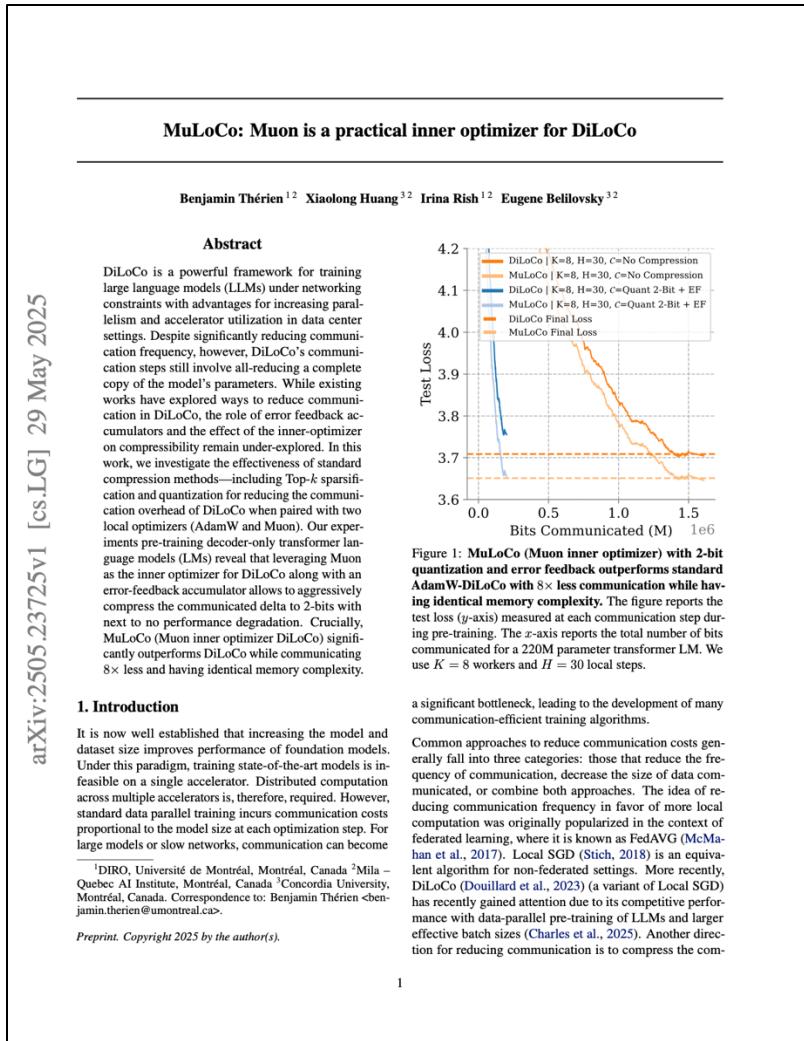
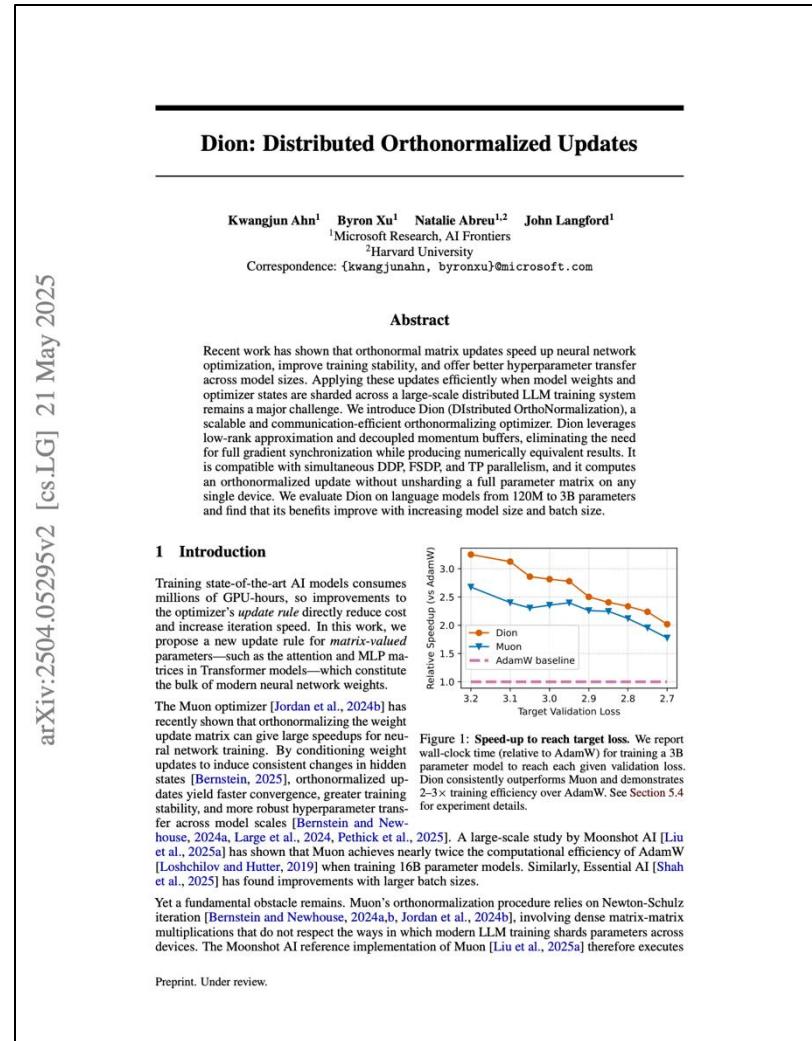
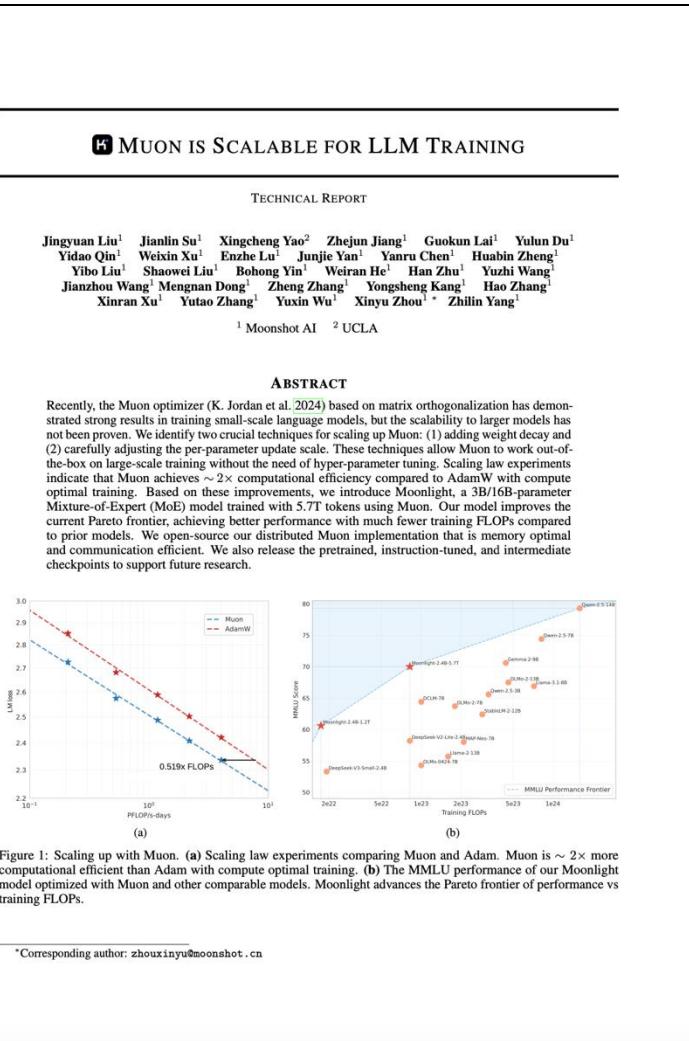
We propose a new way for adapting the learning rate which we found to be the best way for understanding the underlying norm constraint in the parameter space. There

is a well-known early work by Carlsson (1976) on the stochastic spectral descent method, which deforms steepest descent in

the use
minimizer (Gupta
track at the
Dahl

SCG)
momen-

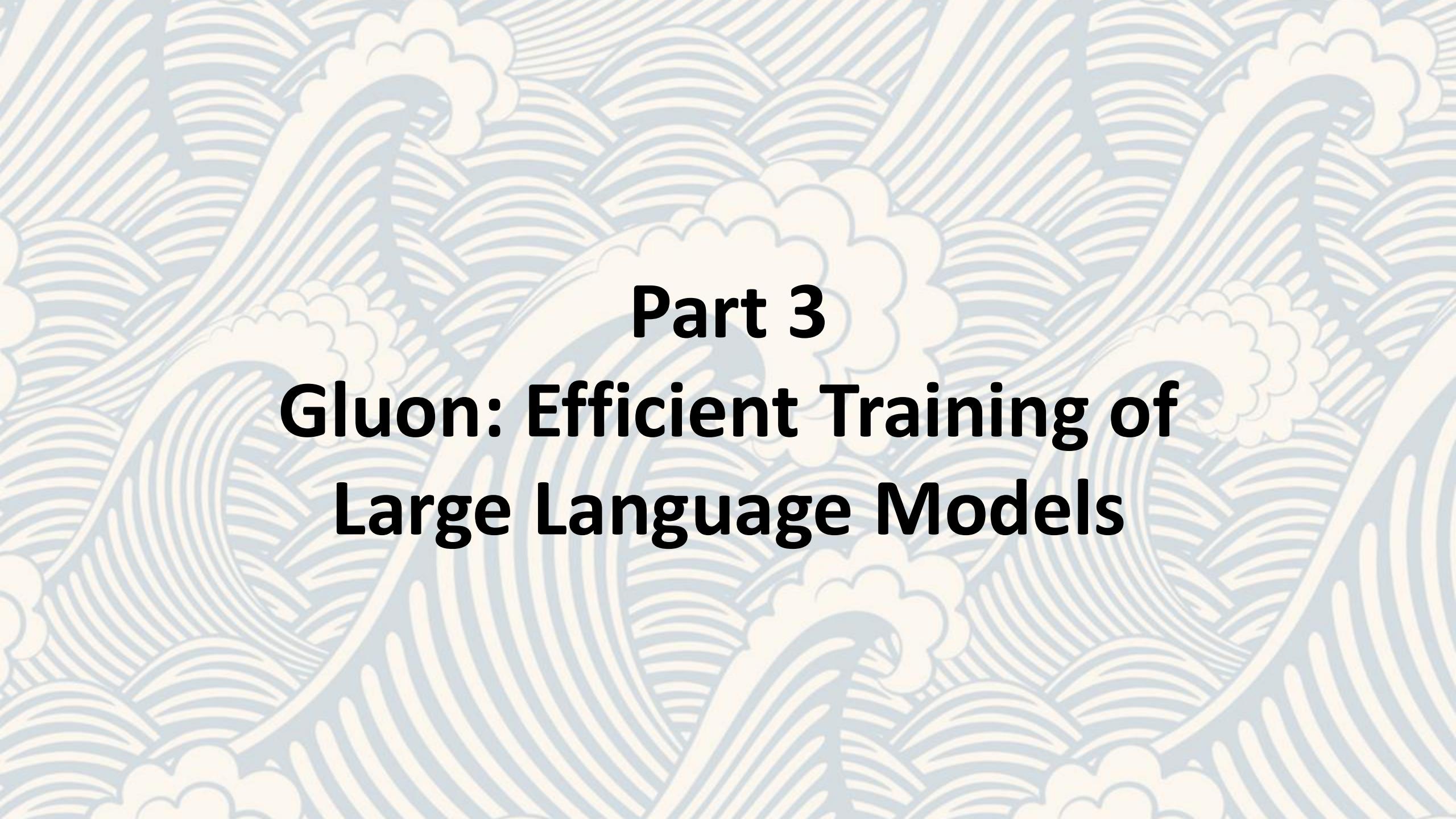
Distributed Muon-Style Optimizers



February 2025

May 2025

May 2025



Part 3

Gluon: Efficient Training of

Large Language Models

Issues with Existing Theory

1. Global vs layer-wise linear minimization

- **Practice:** linear minimization over non-Euclidean balls is performed layer-by layer
- **Theory:** linear minimization over all parameters

2. Learning rates (stepsizes)

- **Practice:** trial & error tuning
- **Theory:** impractically tiny stepsizes

Gluon: Making Muon & Scion Great Again!

(Bridging Theory and Practice of LMO-based Optimizers for LLMs)

Artem Riabinin Egor Shulgin Kaja Gruntkowska Peter Richtárik

King Abdullah University of Science and Technology (KAUST)
Thuwal, Saudi Arabia

Abstract

Recent developments in deep learning optimization have brought about radically new algorithms based on the Linear Minimization Oracle (LMO) framework, such as Muon [12] and Scion [24]. After over a decade of Adam's dominance, these LMO-based methods are emerging as viable replacements, offering several practical advantages such as improved memory efficiency, better hyperparameter transferability, and most importantly, superior empirical performance on large-scale tasks, including LLM training. However, a significant gap remains between their practical use and our current theoretical understanding: prior analyses (1) overlook the layer-wise LMO application of these optimizers in practice, and (2) rely on an unrealistic smoothness assumption, leading to impractically small stepsizes. To address both, we propose a new LMO-based method called Gluon, capturing prior theoretically analyzed methods as special cases, and introduce a new refined generalized smoothness model that captures the layer-wise geometry of neural networks, matches the layer-wise practical implementation of Muon and Scion, and leads to convergence guarantees with strong practical predictive power. Unlike prior results, our theoretical stepsizes closely match the fine-tuned values reported by Pethick et al. [24]. Our experiments with NanoGPT and CNN confirm that our assumption holds along the optimization trajectory, ultimately closing the gap between theory and practice.

1 Introduction

The success of deep learning models across a wide range of challenging domains is inseparable from the optimization algorithms used to train them. As neural networks have grown deeper and datasets larger, optimization has quietly become one of the most consequential components of modern machine learning (ML). Nowhere is this more evident than in the training of large language models (LLMs), which routinely consume thousands of GPU-hours. Adam [15] (and lately AdamW [22])—being effective, relatively reliable, and widely adopted—has for over a decade served as the default choice for this task. While this reliance has powered much of deep learning's progress, it has also exposed the shortcomings of adaptive moment estimation as a one-size-fits-all solution—namely, sensitivity to learning rate schedules, heavy tuning requirements [28], and poor generalization when not carefully calibrated [31]. However, a shift may now be underway. Recent optimizers, such as Muon [12] and Scion [24], represent a significant departure from Adam-type methods: they forgo the adaptive moment estimation in favor of a geometry-aware approach inspired by Frank-Wolfe

Gluon: The Setup

$$\min_{X \in \mathcal{S}} \{ f(X) := \mathbb{E}_{\xi \sim \mathcal{D}} [f_\xi(X)] \}$$

$$\mathcal{S} := \mathcal{S}_1 \otimes \cdots \otimes \mathcal{S}_p$$

$$X = [X_1, \dots, X_p]$$

$$X_i \in \mathcal{S}_i := \mathbb{R}^{n_i \times m_i}$$

Gluon: The Algorithm

$$\min_{X \in \mathcal{S}} \{f(X) := \mathbb{E}_{\xi \sim \mathcal{D}} [f_\xi(X)]\}$$

Stochastic gradient + momentum

$$M_i^k = \beta^k M_i^{k-1} + (1 - \beta^k) \nabla_i f_{\xi^k}(X^k)$$

For all layers $i = 1, 2, \dots, p$ do:

$$X_i^{k+1} = \operatorname{argmin}_{X_i \in \mathcal{S}_i} \left\{ \langle M_i^k, X_i \rangle_{(i)} \mid \|X_i - X_i^k\|_{(i)} \leq t_i^k \right\}$$

$$X^{k+1} = [X_1^{k+1}, \dots, X_i^{k+1}, \dots, X_p^{k+1}]$$

Trainable parameters belonging to layer i

$$X_i \in \mathcal{S}_i := \mathbb{R}^{n_i \times m_i}$$

Trace inner product on $\mathcal{S}_i := \mathbb{R}^{m_i \times n_i}$

$$\langle X_i, Y_i \rangle := \operatorname{Tr} (X_i^\top Y_i)$$

Stepsize/radius for layer i

Non-Euclidean norm on $\mathcal{S}_i := \mathbb{R}^{m_i \times n_i}$
Spectral norm for $i = 1, \dots, p-1$

New Smoothness Assumption

New Assumption: Layer-wise L0-L1 Smoothness

$$\min_{X \in \mathcal{S}} \{f(X) := \mathbb{E}_{\xi \sim \mathcal{D}} [f_\xi(X)]\}$$

Dual norm for layer i

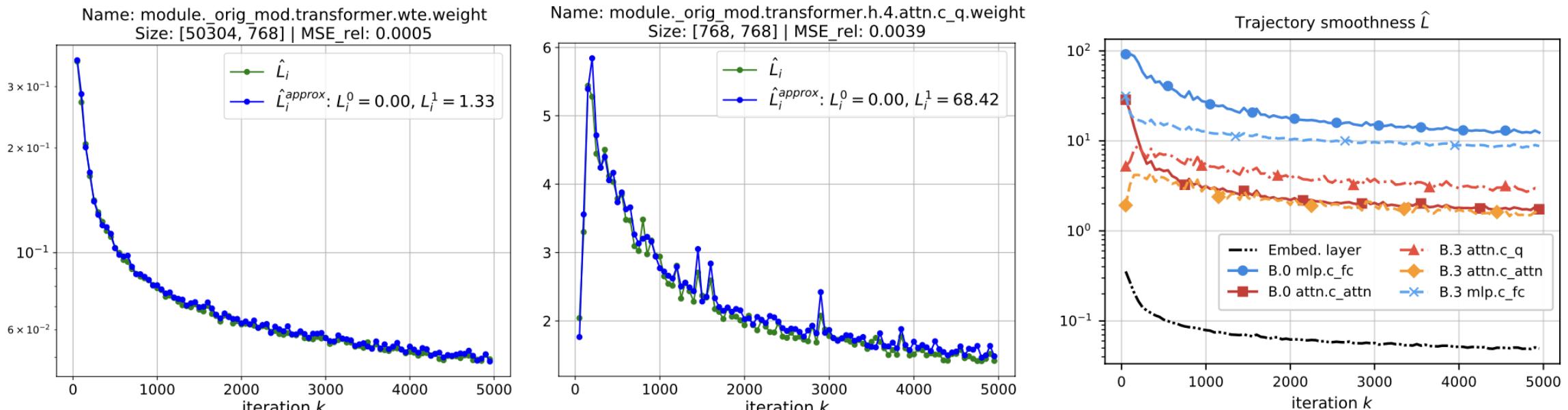
$$\|G_i\|_{(i)*} := \sup_{\|X_i\|_{(i)} \leq 1} \langle G_i, X_i \rangle_{(i)}$$

$$\frac{\|\nabla_i f(X) - \nabla_i f(Y)\|_{(i)*}}{\|X_i - Y_i\|_{(i)}} \leq L_i^0 + L_i^1 \|\nabla_i f(X)\|_{(i)*}, \quad \forall X, Y \in \mathcal{S}, \quad \forall i \in \{1, \dots, p\}$$



Two smoothness constants for each layer i

Experimental Validation of Layer-wise L0-L1 Smoothness



(a) Token embedding matrix from the first/last layer.
(b) Self-attention query matrix from the 4th transformer block.
(c) Trajectory smoothness across different blocks ($B.i$) and layers.

Figure 1: Training NanoGPT on FineWeb validates our layer-wise (L^0, L^1) -smoothness model.

Gluon Convergence Theory

Gluon: Deterministic Variant

$$\min_{X \in \mathcal{S}} \{f(X) := \mathbb{E}_{\xi \sim \mathcal{D}} [f_\xi(X)]\}$$

~~Stochastic gradient + momentum~~

$$M_i^k = \beta^k M_i^{k-1} + (1 - \beta^k) \nabla_{\xi^k} f_{\xi^k}(X^k)$$

For all layers $i = 1, 2, \dots, p$ do:

$$X_i^{k+1} = \operatorname{argmin}_{X_i \in \mathcal{S}_i} \left\{ \langle M_i^k, X_i \rangle_{(i)} \mid \|X_i - X_i^k\|_{(i)} \leq t_i^k \right\}$$

$$X^{k+1} = [X_1^{k+1}, \dots, X_i^{k+1}, \dots, X_p^{k+1}]$$

Trainable parameters belonging to layer i

$$X_i \in \mathcal{S}_i := \mathbb{R}^{n_i \times m_i}$$

Trace inner product on $\mathcal{S}_i := \mathbb{R}^{m_i \times n_i}$

$$\langle X_i, Y_i \rangle := \operatorname{Tr} (X_i^\top Y_i)$$

Stepsize/radius for layer i

Non-Euclidean norm on $\mathcal{S}_i := \mathbb{R}^{m_i \times n_i}$

Spectral norm for $i = 1, \dots, p-1$

Gluon: Deterministic Variant

$$\min_{X \in \mathcal{S}} \{f(X) := \mathbb{E}_{\xi \sim \mathcal{D}} [f_\xi(X)]\}$$

Gradient

$$M_i^k = \nabla f_i(X^k)$$

Trace inner product on $\mathcal{S}_i := \mathbb{R}^{m_i \times n_i}$

$$\langle X_i, Y_i \rangle := \text{Tr}(X_i^\top Y_i)$$

Stepsize/radius for layer i

For all layers $i = 1, 2, \dots, p$ do:

$$X_i^{k+1} = \operatorname{argmin}_{X_i \in \mathcal{S}_i} \left\{ \langle M_i^k, X_i \rangle_{(i)} \mid \|X_i - X_i^k\|_{(i)} \leq t_i^k \right\}$$

$$X^{k+1} = [X_1^{k+1}, \dots, X_i^{k+1}, \dots, X_p^{k+1}]$$

Trainable parameters belonging to layer i

$$X_i \in \mathcal{S}_i := \mathbb{R}^{n_i \times m_i}$$

Non-Euclidean norm on $\mathcal{S}_i := \mathbb{R}^{m_i \times n_i}$
Spectral norm for $i = 1, \dots, p-1$

Gluon Theory 1

$$M_i^k = \nabla f_i(X^k) \quad t_i^k = \frac{\|\nabla_i f(X^k)\|_{(i)*}}{L_i^0 + L_i^1 \|\nabla_i f(X^k)\|_{(i)*}}$$

$$\Delta^0 = f(X^0) - \inf f$$

K

Total number of iterations

$$\min_{k=0, \dots, K-1} \sum_{i=1}^p \|\nabla_i f(X^k)\|_{(i)*} \leq \varepsilon$$

$$\text{MAX} = \max\{L_1^1, \dots, L_p^1\}$$

(maximum of the smoothness constants $\{L_i^1\}_{i=1}^p$)

$$K \geq 2\Delta^0 \cdot \left(\frac{\sum_{i=1}^p L_i^0}{\varepsilon^2} + \frac{\text{MAX}}{\varepsilon} \right)$$

Gluon Theory 2

$$M_i^k = \nabla f_i(X^k) \quad t_i^k = \frac{\|\nabla_i f(X^k)\|_{(i)*}}{L_i^0 + L_i^1 \|\nabla_i f(X^k)\|_{(i)*}}$$

$$\Delta^0 = f(X^0) - \inf f$$

$$K \geq 2\Delta^0$$

Total number of iterations

$$\min_{k=0, \dots, K-1} \sum_{i=1}^p \frac{\text{HM}}{L_i^1} \|\nabla_i f(X^k)\|_{(i)*} \leq \varepsilon$$

$$\text{HM} = \frac{1}{\frac{1}{p} \sum_{i=1}^p \frac{1}{L_i^1}}$$

(harmonic mean of the smoothness constants $\{L_i^1\}_{i=1}^p$)

$$\left(\frac{\sum_{i=1}^p \left(\frac{\text{HM}}{L_i^1} \right)^2 L_i^0}{\varepsilon^2} + \frac{\text{HM}}{\varepsilon} \right)$$



Gluon Theory 3

$$\mathbb{E}_{\xi \sim \mathcal{D}} \left[\|\nabla_i f_\xi(X) - \nabla_i f(X)\|_{(i)*}^2 \right] \leq \sigma^2$$

$$\forall X \in \mathcal{S}, \forall i \in \{1, \dots, p\}$$

$$\Delta^0 = f(X^0) - \inf f$$

$$\min_{k=0, \dots, K-1} \sum_{i=1}^p \frac{1}{12L_i^1} \mathbb{E} \left[\|\nabla_i f(X^k)\|_{(i)*} \right] \lesssim \frac{\Delta^0}{K^{1/4}} + \frac{\sum_{i=1}^p \left(\frac{L_i^0}{(L_i^1)^2} + \frac{\sigma}{L_i^1} \right)}{K^{1/4}}$$

Total number of iterations

Gluon Theory vs Previous Results

Table 1: Comparison of convergence guarantees for Gluon (Algorithms 1 and 2) to achieve $\min_{k=0,\dots,K-1} \sum_{i=1}^p \mathbb{E}[\|\nabla_i f(X^k)\|_{(i)\star}] \leq \varepsilon$, where the $\mathcal{O}(\cdot)$ notation hides logarithmic factors. Notation: K = total number of iterations, (L^0, L^1) = the result holds under layer-wise (L^0, L^1) -smoothness, t_i^k = radius/stepsizes, $1 - \beta^k$ = momentum.

Result	Stochastic?	(L^0, L^1)	Rate	Stepsizes t_i^k	$1 - \beta^k$
[16, Theorem 1]	✗	✗	$\mathcal{O}\left(\frac{1}{K^{1/2}}\right)$	const $\propto \frac{1}{K^{1/2}}$ (b)	—
[16, Theorem 2]	✓	✗	$\mathcal{O}\left(\frac{1}{K^{1/4}}\right)$	const $\propto \frac{1}{K^{3/4}}$ (b)	const $\propto \frac{1}{K^{1/2}}$
[18, Theorem 2.1] ^(a)	✓	✗	$\mathcal{O}\left(\frac{1}{K^{1/4}}\right)$	const $\propto \frac{1}{K^{3/4}}$ (b)	const $\propto \frac{1}{K^{1/2}}$
[24, Lemma 5.4]	✓	✗	$\mathcal{O}\left(\frac{1}{K^{1/4}}\right)$	const $\propto \frac{1}{K^{3/4}}$ (b)	$\propto \frac{1}{k^{1/2}}$
NEW: Theorem 1	✗	✓	$\mathcal{O}\left(\frac{1}{K^{1/2}}\right)$	Adaptive	—
NEW: Theorem 2	✓	✓	$\mathcal{O}\left(\frac{1}{K^{1/4}}\right)$	$\propto \frac{1}{k^{3/4}}$	$\propto \frac{1}{k^{1/2}}$

(a) Applies only to the Muon/Scion update in (14) with $p = 1$.

(b) These stepsizes are impractically tiny since they have an inverse dependence on the total number of iterations K .



The End