
Enhancing Optimizer Stability: Momentum Adaptation of The NGN Step-size

Rustem Islamov¹, Niccolò Ajroldi², Antonio Orvieto^{2,3,4}, and Aurelien Lucchi¹

¹University of Basel, Switzerland

²Max Planck Institute for Intelligent Systems, Germany

³ELLIS Institute Tübingen, Germany

⁴Tübingen AI Center, Germany

Abstract

Modern optimization algorithms that incorporate momentum and adaptive step-size offer improved performance in numerous challenging deep learning tasks. However, their effectiveness is often highly sensitive to the choice of hyperparameters, especially the step-size. Tuning these parameters is often difficult, resource-intensive, and time-consuming. Therefore, recent efforts have been directed toward enhancing the stability of optimizers across a wide range of hyperparameter choices [Schaipp et al., 2024]. In this paper, we introduce an algorithm that matches the performance of state-of-the-art optimizers while improving stability to the choice of the step-size hyperparameter through a novel adaptation of the NGN step-size method [Orvieto and Xiao, 2024]. Specifically, we propose a momentum-based version (NGN-M) that attains the standard convergence rate of $\mathcal{O}(1/\sqrt{K})$ under less restrictive assumptions, without the need for interpolation condition or assumptions of bounded stochastic gradients or iterates, in contrast to previous approaches. Additionally, we empirically demonstrate that the combination of the NGN step-size with momentum results in enhanced robustness to the choice of the step-size hyperparameter while delivering performance that is comparable to or surpasses other state-of-the-art optimizers.

1 Introduction

Adaptive methods such as Adam [Kingma and Ba, 2015] and RMSprop [Hinton et al., 2012] are widely used in machine learning due to their established advantages over (momentum) SGD, particularly in tasks such as training Transformers [Brown, 2020, Touvron et al., 2021, 2023]. These methods adaptively scale the step-size across different dimensions (parameters) based on their respective statistics, effectively acting as a diagonal preconditioning.

Although these methods perform well in practice, existing theoretical analyses typically require stringent assumptions on the noise structure of the stochastic gradients, such as sub-Gaussian noise [Li et al., 2024] or affine noise models [Wang et al., 2024, Zhang et al., 2024a]: Relaxing these assumptions remains an open challenge. Another well-known issue of Adam is its performance sensitivity to the step-size hyperparameter [Wilson et al., 2017, Choi et al., 2019], particularly when training Transformers, where loss spikes are commonly observed [Molybog et al., 2023, Wortsman et al., 2023]. This often necessitates careful adjustments of the hyperparameters throughout the training process [Zhang et al., 2022, Chowdhery et al., 2023], which can be costly in terms of computational resources [Or et al., 2020]. Consequently, there has been growing interest in developing optimization methods that are more robust to hyperparameter selection [Schaipp et al., 2024]. In addition to adapting the step-size, Adam and other state-of-the-art optimizers also rely on momentum [Polyak, 1964], a broadly used technique that has been shown to enhance performance both

theoretically [Cutkosky and Mehta, 2020, Fatkhullin et al., 2024, Islamov et al., 2024b, 2025] and practically [Choi et al., 2019, Fu et al., 2023, Jelassi and Li, 2022]. Besides speeding up convergence, momentum is known as a technique to reduce the variance of stochastic algorithms [Ma and Yarats, 2018, Cutkosky and Orabona, 2019], improving stability as well as generalization in some settings [Jelassi and Li, 2022].

In this work, we address the aforementioned drawbacks of Adam by developing a new algorithm based on the recently proposed NGN step-size [Orvieto and Xiao, 2024], an improved variant of the Stochastic Polyak Step-size [Loizou et al., 2021] that has demonstrated strong resilience to step-size hyperparameter tuning. In particular, NGN was shown never to diverge for any choice of the step-size hyperparameter in the convex setting, and to exhibit strong curvature adaptation properties strengthened by theoretical guarantees. However, the step-size of Orvieto and Xiao [2024] simply adapts the learning rate through a scalar multiplier, leaving to future work the incorporation of momentum and coordinate-wise variants – needed in complex problems such as optimizing transformers, as motivated above. Here, we develop a momentum and step-size adaptive version of NGN designed to enhance robustness in terms of hyperparameter selection. We also present a theoretical analysis alongside a practical evaluation of this approach, showcasing its improvements over current state-of-the-art methods.

In summary, our contributions are as follows:

1. We introduce a new algorithm named NGN-M that combines the NGN step-size with momentum. We theoretically show that NGN-M achieves a convergence rate $\mathcal{O}(1/\sqrt{K})$ in the convex regime without the typical requirements of interpolation or bounded gradient assumptions found in earlier works;
2. We focus on the problem of adapting the step-size rule towards a coordinate-wise diagonal preconditioning. By integrating this diagonal step-size strategy with momentum, we develop a new variant of NGN, called NGN-MD;
3. The theoretical results are supported by extensive empirical validation in various deep learning settings where we demonstrate that NGN-M and NGN-MD not only preserve the robustness property of the NGN step-size, but improve it further in many cases. The step-size hyperparameter resilience comes together with better or comparable performance to state-of-the-art algorithms.

2 Related Works

Polyak Step-size. When training a deep network with standard optimizers, tuning the learning rate is crucial but time-consuming and resource-intensive [Goodfellow et al., 2016]. This issue is at the root of recent research focusing on transferring hyperparameters across architectures at different scales, therefore avoiding expensive tuning pipelines [Yang et al., 2022, 2023, Bordelon et al., 2023]. Yet, in the convex setting, choosing the learning rate can already be difficult – an issue that was studied already in Polyak [1987] and gave rise to the first adaptive method: the Polyak Stepsize (PS). Recently, there has been a renewed interest in adapting PS to modern settings [Loizou et al., 2021, Orvieto et al., 2022, Jiang and Stich, 2024], delivering a theoretically principled way to scale the gradient magnitude during training adaptively. PS-inspired methods have gained increasing interest for their simplicity and adaptability, as they utilize local curvature and smoothness information to accelerate algorithms and facilitate faster convergence. Orvieto and Xiao [2024] recently introduced a variant of the Stochastic Polyak step-size, called NGN, which further enhances the robustness of the step-size hyperparameter and solidifies the link to Gauss-Newton preconditioning. The theoretical analysis in Orvieto and Xiao [2024] demonstrated that NGN does not diverge regardless of the choice of the step-size hyperparameter, and converges fast when the step-size is appropriately tuned. In contrast, the current theory of the SPS step-size with fixed step-size hyperparameters [Loizou et al., 2021] proves convergence to the exact solution only if the interpolation condition holds¹.

¹In our notation, this means that $\sigma_{\text{int}}^2 = 0$.

Table 1: Summary of existing methods exploiting Polyak-type adaptive step-sizes and their convergence guarantees. **Mom.**=Supports momentum; **Diag.**=Supports diagonal step-sizes. σ_{int}^2 is defined in Section 4. x^* defines an optimal solution to equation 4. \mathcal{O} notation hides absolute and problem-dependent constant factors and logarithmic terms in the rate.

Method	Rate ^(a)	Mom.	Diag.	Comments
SPS _{max} [Loizou et al., 2021]	$\mathcal{O}(1/K + \sigma_{\text{int}}^2)$	✗	✗	Conv. to non-vanishing neighbourhood
ALR-SMAG [Wang et al., 2023]	$\mathcal{O}((1-\rho)^K + \sigma_{\text{int}}^2)$	✓	✗	Strong convexity Conv. to non-vanishing neighbourhood
Momo [Schaipp et al., 2024]	$\mathcal{O}(1/\sqrt{K})$	✓	✗	Bounded stoch. gradients Interpolation
Momo-Adam [Schaipp et al., 2024]	✗	✓	✓	Momo framework for Adam
MomSPS _{max} [Oikonomou and Loizou, 2024]	$\mathcal{O}(1/K + \sigma_{\text{int}}^2)$	✓	✗	Conv. to non-vanishing neighbourhood
NGN [Orvieto and Xiao, 2024]	$\mathcal{O}(1/\sqrt{K})$	✗	✗	—
NGN-M (Alg. 1) [This work]	$\mathcal{O}(1/\sqrt{K})$	✓	✗	—
NGN-MDv1 (Alg. 2) [This work]	✗	✓	✓	Combination of NGN-M and RMSprop
NGN-MDv2 (Alg. 2) [This work]	✗	✓	✓	Combination of NGN-M and NGN-D
NGN-D (Alg. 3) [This work]	$\mathcal{O}(1/\sqrt{K})$	✗	✓	—

Polyak Step-size and Heavy-ball Momentum. Heavy-ball momentum methods, stemming from the work of Polyak [1964], have gained significant attention over the years due to their benefits, including acceleration on convex quadratics [Jain et al., 2018, Lee et al., 2022, Bollapragada et al., 2022], convex-like [Wang et al., 2022], and non-convex problems [Cutkosky and Mehta, 2020], as well as their variance reduction abilities [Ma and Yarats, 2018, Cutkosky and Orabona, 2019]. This has led to growing interest in the combination of Polyak step-size and heavy-ball momentum, which is an active area of research [Barré et al., 2020, Saab et al., 2022, Barré et al., 2020, Wang et al., 2023, Oikonomou and Loizou, 2024, Gower et al., 2025]. Recently, Schaipp et al. [2024] demonstrated that a geometrically principled combination of SPS and momentum leads to lower sensitivity to the step-size hyperparameter, although they did not provide strong theoretical convergence guarantees.

Diagonal Polyak Step-size. Coordinate-wise adaptive step-sizes are essential in training Transformer architectures due to the varying parameter-wise scaling and conditioning of the problem [Noci et al., 2022, Zhang et al., 2024b]. Algorithms employing diagonal step-sizes, such as Adam and SignSGD [Bernstein et al., 2018], typically outperform non-diagonal methods in language modeling tasks by addressing issues such as class imbalance (where certain words appear more frequently than others) [Kunstner et al., 2023, 2024] and heavy-tailed noise [Zhang et al., 2019, 2020, Compagnoni et al., 2025]. It is, therefore, paramount in current setups to deliver adaptive step-size improvements targeted to the coordinate-wise (diagonal) regime. However, most Polyak-step-size-based algorithms only focus on a single step-size for all parameters [Loizou et al., 2021, Wang et al., 2023, Gower et al., 2021, Oikonomou and Loizou, 2024, Orvieto and Xiao, 2024]. Only a few works propose a diagonal-wise modification of Polyak-step-size by either using Adam preconditioner [Schaipp et al., 2024] as a weight matrix or incorporating second-order information from the objective function [Li et al., 2022, Richtárik et al., 2024].

Table 1 provides a theoretical comparison of various Polyak step-size-based algorithms that incorporate momentum and/or diagonal step-size, highlighting the differences between the theoretical results presented in this work and those from prior works.

3 Algorithm design of NGN-M and NGN-D

In [Orvieto and Xiao \[2024\]](#), the NGN step-size is derived by applying a Gauss–Newton update on a regularized first-order expansion of $r(x) := \sqrt{f(x)}$. At the current point x^k , they linearized $r(x^k + p) \approx r(x^k) + \nabla r(x^k)^\top p$. Thus the next iterate is given as $x^{k+1} = x^k + p^k$ where

$$p^k := \arg \min_p [(r(x^k) + \nabla r(x^k)^\top p)^2 + \frac{1}{2c} \|p\|^2]. \quad (1)$$

It turns out that the problem above has a closed-form solution

$$p^k = -\gamma_k \nabla f(x^k) \text{ where } \gamma_k := \frac{c}{1 + \frac{c}{2f(x^k)} \|\nabla f(x^k)\|^2},$$

with γ_k representing the NGN step-size. In [Orvieto and Xiao \[2024\]](#), convergence guarantees were established for both convex and general non-convex settings. Importantly, the convex analysis shows that NGN exhibits a non-divergence property, regardless of the step-size hyperparameter c (see Theorem 4.5 in [\[Orvieto and Xiao, 2024\]](#)). Due to this property, the NGN step-size is a strong candidate to achieve better robustness w.r.t. the choice of the step-size.

3.1 How to Add Momentum and What to Expect?

There are several approaches to combining the adaptive Polyak-type step-size with heavy-ball momentum. Broadly, existing algorithms can be divided into two categories: the first category involves computing the Polyak step-size in the usual manner and incorporating it into the standard heavy-ball update [\[Oikonomou and Loizou, 2024\]](#). In contrast, algorithms from the second category first determine an update direction using exponential weighted averaging of the stochastic gradient and momentum variable, and then compute the Polyak-type step-size based on the computed direction [\[Wang et al., 2023, Schaipp et al., 2024\]](#). Following this principled approach, we test two possible versions for combining the NGN step-size and momentum:

$$\text{Ver.1 : } \begin{cases} \gamma_k = \frac{c}{1 + \frac{c}{2f_{S_k}(x^k)} \|\nabla f_{S_k}(x^k)\|^2} \\ m^k = \beta m^{k-1} + (1 - \beta) \gamma_k \nabla f_{S_k}(x^k) \\ x^{k+1} = x^k - m^k \end{cases} \quad \text{Ver.2 : } \begin{cases} m^k = \beta m^{k-1} + (1 - \beta) \nabla f_{S_k}(x^k) \\ \gamma_k = \frac{c}{1 + \frac{c}{2f_{S_k}(x^k)} \|m^k\|^2} \\ x^{k+1} = x^k - \gamma_k m^k \end{cases}.$$

Before we proceed, we should answer the question: “*What do we expect from the combination of NGN step-size and momentum?*”. First, we aim to preserve, and ideally enhance, NGN’s robustness to the step-size hyperparameter. Additionally, we seek improved performance, achieving accelerated convergence akin to the advantage of SGD with momentum (SGDM) over standard SGD in convex settings. With these goals in mind, we now show that version 1 meets all of these criteria, while version 2 is less suitable. To gain some intuition regarding the performance of these two variants, we start by conducting a simple experiment on a quadratic function $f(x) = \frac{1}{2} \|Ax - b\|^2$ where A is a data matrix from the normalized Diabetes dataset [\[Smith et al., 1988\]](#) and b is a vector of labels. Based on the results from Figure 1 (left), we observe that variant 1 achieves accelerated convergence as SGDM for middle-range step-size hyperparameters ($c \in \{10^1, 10^2\}$) and does not diverge for large step-size parameter ($c \in \{10^3\}$). Conversely, version 2 has a worse convergence rate than version 1 for middle-range step-size parameters and diverges for large ones. Therefore, we theoretically analyze and practically test version 1, which we call NGN-M.

3.2 Evidence of Robustness of NGN-M

To illustrate the advantages of the design choice NGN-M, we first consider the Rosenbrock function $f(x, y) = (x - 1)^2 + 100(y - x^2)^2$, whose minimizer is at $(1, 1)$. Starting from $(-1.2, 1)$, we run both NGN-M and SGDM over a wide range of constant step-size hyperparameters $\{10^{-3}, \dots, 10^2\}$. As shown in Figure 1, we observe that (i) for small step-size hyperparameter both methods successfully

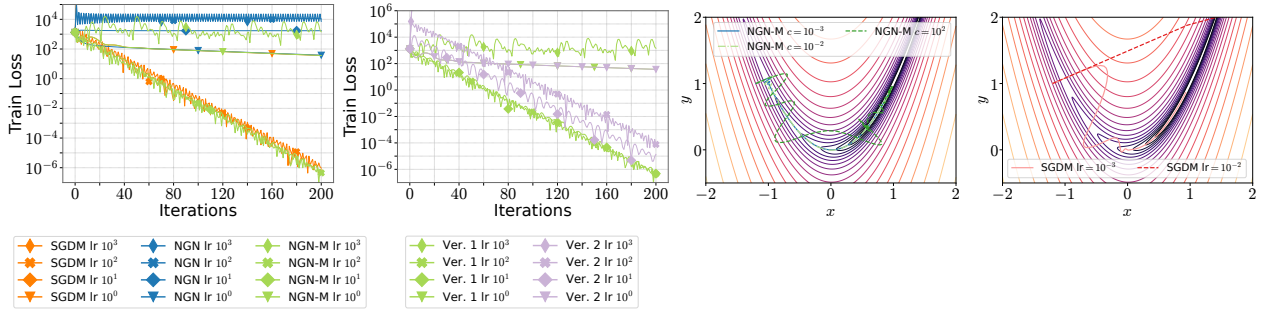


Figure 1: **Left:** Comparison of SGDM, NGN, NGN-M for linear regression on normalized Diabetes dataset varying a step-size hyperparameter. **Second left:** Comparison of two options on how momentum can be used in combination with NGN step-size. **Third and fourth:** Comparison of SGDM and NGN-M on the Rosenbrock function.

converge to $(1, 1)$; (ii) SGDM already diverges for the step-size hyperparameter 10^{-2} ; By contrast, NGN-M remains stable even up to $c = 10^2$, thanks to its adaptive step-size that automatically adjusts with the local curvature. Figure I.3 further traces the optimization trajectories: NGN-M converges reliably for every tested value of c , whereas SGDM fails outside its narrow stability window. Finally, in Appendix I.1 we repeat these experiments on a synthetic multimodal function and find that NGN-M consistently finds the global minimum, while SGDM typically becomes trapped in a nearby suboptimal local minimum.

3.3 Diagonal Step-size for NGN

We propose two alternatives to make NGN step-size coordinate-wise adaptive. In the first approach, we modify an approach of (1): The next iterate x^{k+1} is obtained by minimizing an approximation of the regularized first-order Taylor expansion of $r(x) := \sqrt{f(x)}$ around x^k , namely, $x^{k+1} = x^k + p^k$ where for a preconditioning matrix Σ_k

$$p^k = \arg \min_p \left[(r(x^k) + \nabla r(x^k)^\top p)^2 + \frac{1}{2c} \|p\|_{\Sigma_k}^2 \right]. \quad (2)$$

The intuition is that $\Sigma_k \in \mathbb{R}^{d \times d}$ can penalize each parameter with its own weight while in vanilla NGN the penalization is the same for all parameters, and f is an objective function we aim to minimize. Performing simple derivations (see Appendix G), we obtain the following update rule

$$x^{k+1} = x^k - \frac{c}{1 + \frac{c}{2f(x^k)} \|\nabla f(x^k)\|_{\Sigma_k}^2} \Sigma_k^{-1} \nabla f(x^k). \quad (3)$$

Note that by choosing Σ_k to be an identity matrix, the step-size γ_k in (3) reduces to the vanilla NGN step-size.

Alternatively, we can adopt a simpler, coordinate-wise rule: For each parameter j , we replace the full gradient norm in the NGN step-size with its own partial derivative $\nabla_j f_{S_k}(x^k)$. Both of the described per-coordinate variants can be further adjusted by an RMSprop-style preconditioner $\mathbf{D}_k = \text{diag}((\mathbf{D}_k)_{(1)}, \dots, (\mathbf{D}_k)_{(d)})$ and lead to the following update rule (see Alg. 2 for a full description)

$$\text{NGN-MDv1} : \begin{cases} \gamma_k = \frac{c}{1 + \frac{c}{2f(x^k)} \|\nabla f_{S_k}(x^k)\|_{\mathbf{D}_k}^2} \\ \Sigma_k^{-1} = \gamma_k \mathbf{D}_k^{-1} \end{cases} \quad \text{NGN-MDv2} : \begin{cases} \gamma_k^{(j)} = \frac{c/(\mathbf{D}_k)_{(j)}}{1 + \frac{c/(\mathbf{D}_k)_{(j)}}{2f(x^k)} (\nabla_j f_{S_k}(x^k))^2} \\ \Sigma_k^{-1} = \text{diag}(\gamma_k^{(1)}, \dots, \gamma_k^{(d)}) \end{cases}$$

$$x^{k+1} = x^k - (1 - \beta_1) \Sigma_k^{-1} \nabla f_{S_k}(x^k) + \beta_1 (x^k - x^{k-1})$$

We highlight that both versions have the same number of hyperparameters as Adam. From an empirical evaluation of two versions of NGN-MD in Figure 2, we observe that the first choice improves the performance of NGN-M while maintaining robustness to step-size hyperparameter. A more detailed discussion on the two versions of NGN-MD algorithms is deferred to Appendix G.1.

Algorithm 1 NGN-M

```

1: Input:  $x^{-1} = x^0 \in \mathbb{R}^d$ , step-size hyperparameter  $c > 0$ , momentum parameter  $\beta \in [0, 1)$ 
2: for  $k = 0, 1, \dots, K - 1$  do
3:   Sample a batch  $S_k \subseteq [n]$ 
4:    $\gamma_k = \frac{c}{1 + \frac{c}{2f_{S_k}(x^k)} \|\nabla f_{S_k}(x^k)\|^2}$ 
5:    $x^{k+1} = x^k - (1 - \beta)\gamma_k \nabla f_{S_k}(x^k) + \beta(x^k - x^{k-1})$ 
6: end for

```

In the special case $\beta_1 = 0$ and $\Sigma_k = \mathbf{I}$, NGN-MDv2 reduces to NGN-D (Alg. 3). To the best of our knowledge, NGN-D is the first algorithm that uses a per-parameter Polyak-type step-size while achieving the standard $\mathcal{O}(1/\sqrt{K})$ rate under smoothness and bounded noise variance assumptions; see detailed discussion in Appendix C.

4 Theoretical Analysis of NGN-M

4.1 Problem Formulation and Notation

We consider the classic Empirical Risk Minimization (ERM) problem that typically appears when training machine learning models, namely,

$$\min_{x \in \mathbb{R}^d} [f(x) := \frac{1}{n} \sum_{i=1}^n f_i(x)], \quad (4)$$

where x are the parameters of a model we aim to train, n is the number of data points in the dataset, d is the number of parameters, x^* denotes the solution to equation 4, and f_i represents the loss associated with the i -th data point/batch. We assume that each f_i is differentiable and non-negative² and that the global optimal value is bounded, i.e. $f^* = \arg \min_x f(x) \in \mathbb{R}$. Moreover, we assume that we have access to mini-batch stochastic losses f_S during training such that $f_S^* := \arg \min_x f_S(x) < \infty$ for any $S \subseteq [n]$ picked uniformly at random.

We analyze the convergence of NGN-M under assumptions that are often used in the analysis of the Polyak step-size [Loizou et al., 2021, Orvieto et al., 2022, Orvieto and Xiao, 2024, Oikonomou and Loizou, 2024, Schaipp et al., 2024].

Assumption 1. Each f_i is convex and L -smooth, i.e., for all $x, y \in \mathbb{R}^d$ and $i \in [n]$ we have $\langle \nabla f_i(x), y - x \rangle \geq f_i(x) - f_i(y)$ and $\|\nabla f_i(x) - \nabla f_i(y)\| \leq L\|x - y\|$.

Assumption 2. The interpolation $\sigma_{\text{int}}^2 := \mathbb{E}_S[f^* - f_S^*]$ and positive $\sigma_{\text{pos}}^2 := \mathbb{E}_S[f_S^*]$ errors are bounded. We say that the interpolation holds if $\sigma_{\text{int}}^2 = 0$, where S is a sampled mini-batch.

4.2 Convergence Guarantees

Theorem 1. Let Assumptions 1, 2 hold. Let the step-size hyperparameter $c > 0$ and the momentum parameter $\beta = \frac{\lambda}{1+\lambda}$ be constants where $\lambda \leq \min\{cL, 0.5(1+cL)^{-1}(1+2cL)^{-1}\}$. Then the iterates of NGN-M (Alg. 1) satisfy

$$\mathbb{E} [f(\bar{x}^{K-1}) - f(x^*)] \leq \frac{\|x^0 - x^*\|^2(1+2cL)^2}{cK} + 8cL(1+2cL)^2\sigma_{\text{int}}^2 + 2cL \max\{2cL - 1, 0\} \sigma_{\text{pos}}^2,$$

where \bar{x}^{K-1} is chosen uniformly at random from $\{x^0, \dots, x^{K-1}\}$. Moreover, if we set $c = \mathcal{O}(1/\sqrt{K})$ then we obtain $\mathbb{E} [f(\bar{x}^{K-1}) - f(x^*)] \leq \mathcal{O}(1/\sqrt{K})$.

The convergence of NGN-M is provided in the convex setting, which is motivated by recent works that observe convex-like structures in the loss landscape of neural networks [Islamov et al., 2024a, Hoang et al., 2024] and agreement between convex theory and practice [Schaipp et al.,

²Common losses, e.g. cross-entropy, satisfy this condition.

Algorithm 2 NGN-MD

```
1: Input:  $x^0 \in \mathbb{R}^d$ , step-size hyperparameter  $c > 0$ , momentum parameters  $\beta_1, \beta_2 \in [0, 1)$ , stabilization parameter  $\varepsilon > 0$ , second-order momentum  $v^0 = 0$ 
2: for  $k = 0, 1, \dots, K - 1$  do
3:   Sample a batch  $S_k \subseteq [n]$ 
4:    $v^k = \beta_2 v^{k-1} + (1 - \beta_2)(\nabla f_{S_k}(x^k) \odot \nabla f_{S_k}(x^k))$ 
5:    $\mathbf{D}_k = \text{diag}(\varepsilon \mathbf{I} + \sqrt{v^k / (1 - \beta_2^k)})$ 
6:   For NGN-MDv1:  $\gamma_k = \frac{c}{1 + \frac{c}{2f_{S_k}(x^k)} \|\nabla f_{S_k}(x^k)\|_{\mathbf{D}_k^{-1}}^2}$ 
7:   For NGN-MDv1:  $\Sigma_k^{-1} = \gamma_k \mathbf{D}_k^{-1}$ 
8:   For NGN-MDv2:  $\Sigma_k^{-1} = \text{diag}(\gamma_k^{(1)}, \dots, \gamma_k^{(d)})$  where  $\gamma_k^{(j)} = \frac{c / (\mathbf{D}_k)_{(j)}}{1 + \frac{c}{2f_{S_k}(x^k) \cdot (\mathbf{D}_k)_{(j)}} (\nabla_j f_{S_k}(x^k))^2}$ 
9:    $x^{k+1} = x^k - (1 - \beta_1) \Sigma_k^{-1} \nabla f_{S_k}(x^k) + \beta_1 (x^k - x^{k-1})$ 
10: end for
```

2025]. Importantly, we show that (i) when the constant c is sufficiently small, NGN-M attains the same convergence rate as SGDM [Garrigos and Gower, 2023]. Moreover, for any choice of c , we demonstrate that the NGN-M iterates provably converge to a neighborhood of the optimum and thereafter remain within it; (ii) Unlike prior works, our analysis does not rely on strong assumptions such as bounded gradients, interpolation, or a bounded domain; (iii) For small values of c , NGN-M converges to the exact solution while algorithms such as MomSPS and ALR-SMAG were shown to converge up to a non-vanishing neighborhood of the solution only³. Notably, the non-vanishing neighborhood disappears when the problem satisfies interpolation: We refer to Table 1 for more details and exact rates; (iv) The momentum parameter β is theoretically recommended to be set sufficiently small. A default value of $\beta = 0.9$ is commonly used and works well in our experiments. This discrepancy between theoretical guidance and practical implementation has also been observed in prior works on momentum [Ghadimi et al., 2015, Liu et al., 2020, Wang et al., 2023, 2022, Oikonomou and Loizou, 2024]. Interestingly, for simple functions we can establish convergence even when β is large (see Appendix F), indicating that the small- β requirement may be an artifact of the existing proving techniques rather than an inherent algorithmic limitation of NGN-M. We leave a comprehensive study of arbitrary β values across general convex objectives for future work; (v) While Theorem 1 requires knowing the total iteration count K to ensure convergence, this might be impractical: We therefore also prove convergence using a diminishing step-size of order $1/\sqrt{k}$ in Appendix E, which removes the need to preset K ; (vi) Finally, we corroborate our analysis as we run NGN-M with the theory-derived values of c to a quadratic problem that satisfies all our assumptions: We observe NGN-M’s rapid convergence with theoretical step-size hyperparameters in practice—see Appendix I.3 and Figure I.4 therein.

Key Ingredients of the Proof. We discuss the key steps of the proof to highlight the main challenges in the analysis.

First, we make use of the Iterative Moving Average (IMA) formulation of momentum [Sebbouh et al., 2021]. Specifically, we define a sequence of virtual iterates $\{z^k\}$ whose update rule is of the form

$$z^{k+1} = x^k - \gamma_k \nabla f_{S_k}(x^k), \quad x^{k+1} = \frac{\lambda}{1+\lambda} x^k + \frac{1}{1+\lambda} z^{k+1}, \quad \text{where } z^0 := x^0 \text{ and } \beta = \frac{\lambda}{1+\lambda}.$$

Next, one of the key technical strategies we follow is splitting the step-size γ_k into two parts: a fixed term $\rho = \frac{c}{(1+cL)(1+2cL)} = \mathcal{O}(c)$ and a changing term $\tilde{\gamma}_k \leq \frac{3c^2L}{1+2cL} = \mathcal{O}(c^2)$. This decomposition of the step-size γ_k enables us to regulate the balance between the descent term, which drives improvement in the objective, and the error term, which reflects possible inaccuracies. More precisely, the descent term is weighted by c while the error term proportional to σ_{int}^2 is weighted by c^2 , which suggests that c has to be chosen to tradeoff the two terms to lead to the exact convergence similarly

³In fact, this is an inherited property of SPS analysis from [Loizou et al., 2021].

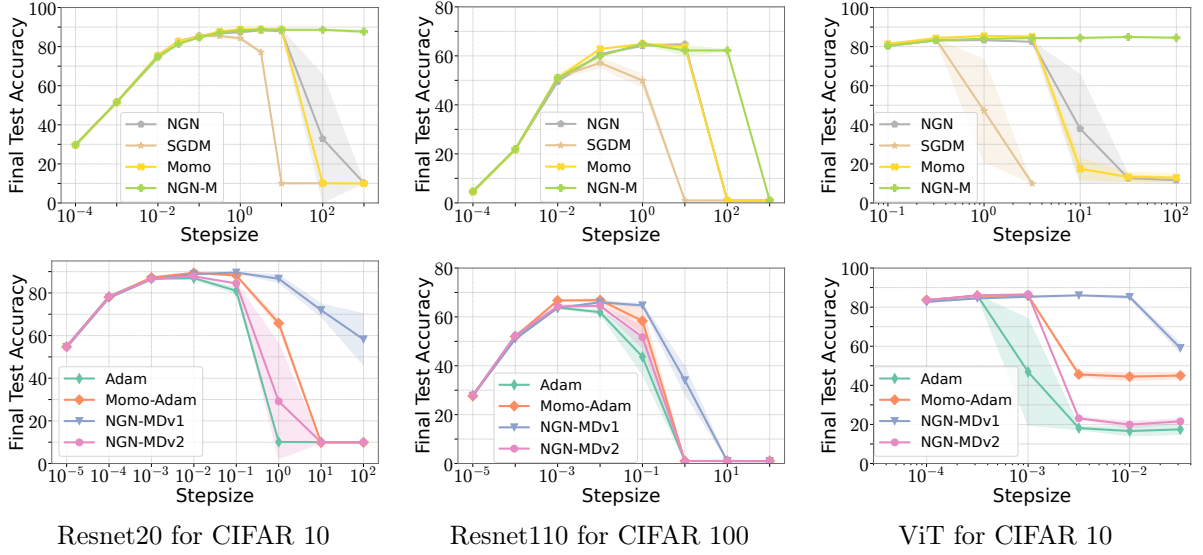


Figure 2: Stability performance of algorithms varying step-size hyperparameter (c for NGN-M, NGN-MDv1 and NGN-MDv2, α_0 for Momo and Momo-Adam, and step-size for SGDM and Adam). For NGN-M and NGN-MDv1, we observe that the range of the step-size hyperparameters that provide competitive performance is wider than that for other algorithms. We refer to Figures J.1 to J.3, J.5 and J.8 for train loss stability and for the results on additional workloads.

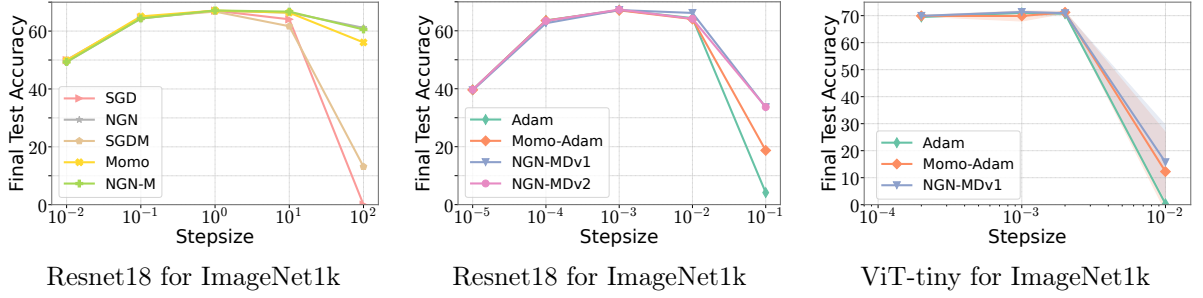


Figure 3: Stability performance on ImageNet1k varying the step-size hyperparameter. NGN-M and NGN-MDv1 achieve higher accuracy for a wider range of the step-size hyperparameters. We refer to Figure J.4 for results on train loss stability and additional results on ImageNet32.

to the standard analysis of SGD [Garrigos and Gower, 2023]. In contrast, MomSPS and Momo algorithms achieve the exact convergence only under the interpolation regime.

5 Experiments

We now turn to the empirical evaluation of the proposed algorithms against several benchmarks. The detailed experiment setup, including the choice of hyperparameters as well as additional experimental results and details, can be found in Appendix J. The best performance of algorithms is reported in Tables 4 (momentum-based algorithms), 5 (algorithms with momentum and component-wise step-size), and 6 (algorithms with component-wise step-size). For clarity and quick reference, all links to the paper’s empirical results are summarized in Table 3.

Comparison on Standard Benchmarks. First, we test the performance of NGN-M against other methods that use momentum, such as SGDM, Momo, MomSPS, ALR-SMAG, and NGN. The tests include the training of Resnet20 [He et al., 2016] and ViT [Dosovitskiy et al., 2021] on the CIFAR10 dataset [Krizhevsky et al., 2009], and Resnet110 on CIFAR100. Second, we test the performance of NGN-MD against Adam and Momo-Adam that – contrary to NGN-M – both use

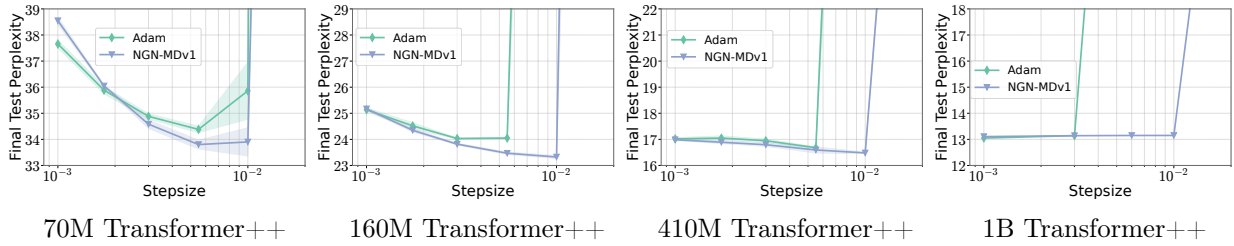


Figure 4: Language Modeling on SlimPajama. Stability comparison with respect to the step-size hyperparameter across different model sizes and optimizers. At all model capacities, NGN-MDv1 achieves similar or lower perplexity, showing better stability and improved performance at larger learning rates. We refer to Figures J.11 to J.14 for the results that report update magnitude when training 160M model and training dynamics across all model sizes.

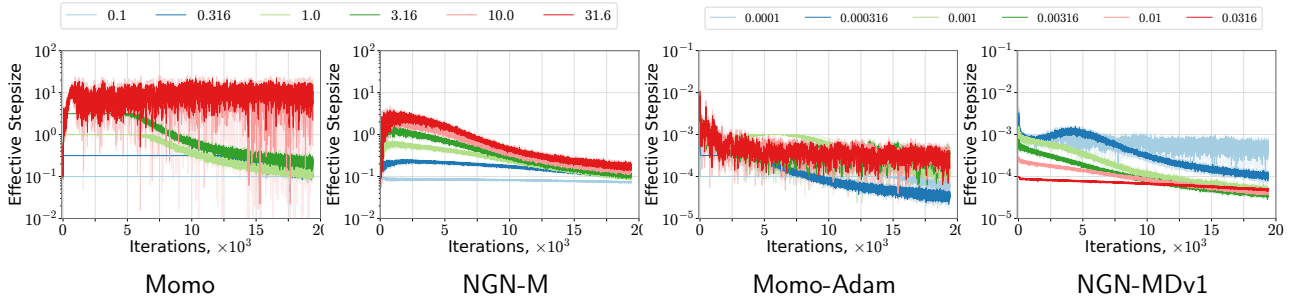


Figure 5: The step-size of Momo, NGN-M (two left), Momo-Adam and NGN-MDv1 (two right) during the training of ViT on CIFAR10. We demonstrate the step-sizes τ_k for Momo and Momo-Adam and γ_k for NGN-M and NGN-MDv1 varying step-size parameters α_0 and c correspondingly. We refer to Figures J.9 and J.10 for the results in training Resnet20.

component-wise preconditioning. All experiments in this section do not use learning rate schedulers or weight decay.

From Tables 4 and 5 we observe that the best performance of NGN-M and NGN-MDv1 matches the results of other algorithms: NGN-M and NGN-MDv1 exhibit competitive performance across all settings we tested. Importantly, NGN-M and NGN-MDv1 demonstrate significantly greater robustness to the choice of the step-size hyperparameter. Indeed, Figure 2 shows that the range of step-size hyperparameter that allows NGN-M and NGN-MDv1 to perform optimally is much wider: We can, for instance, use step-sizes that are 1-2 orders of magnitude larger than the optimal one without a significant drop in the performance. This is particularly evident when training ResNet20 and ViT models. Besides, we clearly observe that momentum consistently improves the stability of NGN across all settings. We refer to Appendix J for additional ablation studies against other optimizers and results when training NLP models.

Vision Experiments on ImageNet. Having observed promising results on workloads of small and medium size, we switch to larger tasks and datasets. We first train a ResNet18 on ImageNet1k [Deng et al., 2009]. This represents the first task in which we pair our proposed algorithms with a learning rate schedule. As illustrated in Figure 3 and Table 4, NGN-M achieves the highest validation accuracy, while exhibiting higher robustness across larger step-sizes, improving over both NGN and Momo. Among adaptive methods, NGN-MDv1 compares favorably against Adam and MomoAdam, while once again achieving higher performance on a wider range of learning rates (Table 5). Appendix J.4 reports additional ablations on ImageNet32 and train loss stability results.

Finally, we test the effectiveness of the proposed algorithms on vision transformers [Dosovitskiy et al., 2021]. These models are trained for a longer horizon compared to convolutional architectures, are notoriously sensitive to initial learning rate, and require adaptive step-sizes. We follow the protocol of Schaiipp et al. [2024], which includes cosine annealing, but without any weight decay

regularization. As highlighted in Figure 3 and Table 5, NGN-MDv1 achieves the highest validation accuracy across adaptive methods. Moreover, at a larger learning rate, Adam diverges, whereas both MomoAdam and NGN-MDv1 maintain more stable training dynamics.

Language Modeling. Pre-training Large Language Models represents a challenging optimization task. To achieve competitive performance, optimizers with adaptive step-size are needed, and preventing instabilities in low-precision training often requires careful hyperparameter tuning.

To evaluate the capability of NGN-MDv1 in this setting, we train decoder-only transformers [Radford et al., 2019] with 70M, 160M, 410M, and 1B parameters around Chinchilla optimum [Hoffmann et al., 2022] on SlimPajama-627B [Soboleva et al., 2023]. For each model, we retune the learning rate, using 3 seeds for the first three models and 1 seed for the 1B. Appendix J provides additional details about the training and tokenization pipeline.

Figure 4 and Table 5 report the final validation perplexity when training language models varying a model size. We note that NGN-MDv1 matches the performance of Adam across all model sizes. However, NGN-MDv1 achieves competitive performance even for a step-size hyperparameter $c = 10^{-2}$ while Adam’s performance drops significantly. This phenomenon is consistent across all scales we tested, suggesting that the optimal learning rate of NGN-MDv1 is shifted towards larger values, but also that the algorithm is less sensitive to such a hyperparameter. We additionally discuss how to introduce weight decay in NGN-MDv1 and report additional ablations on its role in this training task in Appendix H.

Effective Step-size of NGN-M and NGN-MDv1. The first observation from the results in Figure 5 is that the effective step-size of NGN-M and NGN-MDv1 is always adaptive: if the step-size hyperparameter c is large enough the effective step-size sharply increases in the beginning up to a peak, and then it gradually decreases till the end of the training. From this perspective, NGN-M and NGN-MDv1 step-sizes are close to annealing step-size schedulers widely used in practice. In contrast, the effective step-size of Momo and Momo-Adam is not adaptive for sufficiently large step-size hyperparameter α_0 during the initial part or all of the training. In other words, these algorithms reduce to SGDM and Adam, which is one of the reasons for the reduced resilience property of Momo and Momo-Adam in comparison with NGN-M and NGN-MDv1. The effective step-sizes in training Resnet20 are provided Figures J.9 and J.10 while comparison against Adam’s effective step-size is reported in Figures J.6 and J.7. Moreover, we report the update magnitudes when training a 160M language model in Figures J.11 to J.13. All aforementioned results demonstrate that the NGN step-size is more conservative: it decreases the effective step-size when necessary to stabilize the training, even for large values of the step-size hyperparameter c . This feature is a key factor behind robustness of NGN-M and NGN-MDv1 in practice.

6 Conclusion and Future Work

This work introduced several novel adaptations of the NGN step-size method, incorporating support for momentum and/or diagonal step-size. We provided a theoretical analysis of the convergence rates for these algorithms and conducted an extensive empirical evaluation of their performance. The experimental results show that combining momentum with the NGN step-size yields high robustness to step-size hyperparameter choices and performs competitively with state-of-the-art algorithms across various settings.

Given the significant complexity of the task, we defer the theoretical explanation of the step-size resilience properties of NGN-M for large values of β and analysis in the non-convex setting to future work. Furthermore, while the two proposed methods for incorporating weight decay into NGN-MDv1 outperform AdamW in training language models, they still exhibit some sensitivity to the step-size hyperparameter. This may, in part, be due to the limited understanding of the expected effects of the weight decay technique, a topic that requires further investigation. We acknowledge that computing NGN step-size at a large scale may cause runtime overhead, and discuss this limitation in Appendix G.2 by providing train and optimization times. We also recognize that integrating

NGN-MDv1 with advanced parallelism schemes—such as Tensor Parallelism [Shoeybi et al., 2019] or ZeRO-2 [Rajbhandari et al., 2020]—introduces additional compute and communication overhead, and will require further adaptation of the algorithm. Nevertheless, our results provide valuable guidance for developing inherently more stable optimizers. As a next step, it would be fascinating to investigate whether the resilience of emerging methods like Muon [Jordan et al., 2024] can be further improved by incorporating the NGN step-size.

Acknowledgement

Rustem Islamov and Aurelien Lucchi acknowledge the financial support of the Swiss National Foundation, SNF grant No 207392. Antonio Orvieto acknowledges the financial support of the Hector Foundation.

References

- Maksym Andriushchenko, Francesco D’Angelo, Aditya Varre, and Nicolas Flammarion. Why do we need weight decay in modern deep learning? *arXiv preprint arXiv:2310.04415*, 2023. (Cited on pages 38 and 41)
- Mathieu Barré, Adrien Taylor, and Alexandre d’Aspremont. Complexity guarantees for polyak steps with momentum. In *Proceedings of Thirty Third Conference on Learning Theory*, 2020. (Cited on page 3)
- Jeremy Bernstein, Yu-Xiang Wang, Kamyar Azizzadenesheli, and Animashree Anandkumar. signsgd: Compressed optimisation for non-convex problems. In *International Conference on Machine Learning*, 2018. (Cited on pages 3 and 19)
- Sid Black, Stella Biderman, Eric Hallahan, Quentin Anthony, Leo Gao, Laurence Golding, Horace He, Connor Leahy, Kyle McDonell, Jason Phang, Michael Pieler, USVSN Sai Prashanth, Shivan-shu Purohit, Laria Reynolds, Jonathan Tow, Ben Wang, and Samuel Weinbach. Gpt-neox-20b: An open-source autoregressive language model. *arXiv preprint arXiv: 2204.06745*, 2022. (Cited on page 45)
- Raghu Bollapragada, Tyler Chen, and Rachel Ward. On the fast convergence of minibatch heavy ball momentum. *arXiv preprint arXiv:2206.07553*, 2022. (Cited on page 3)
- Blake Bordelon, Lorenzo Noci, Mufan Bill Li, Boris Hanin, and Cengiz Pehlevan. Depthwise hyperparameter transfer in residual networks: Dynamics and scaling limit. *arXiv preprint arXiv:2309.16620*, 2023. (Cited on page 2)
- Tom B Brown. Language models are few-shot learners. *arXiv preprint arXiv:2005.14165*, 2020. (Cited on page 1)
- Xiangning Chen, Chen Liang, Da Huang, Esteban Real, Kaiyuan Wang, Hieu Pham, Xuanyi Dong, Thang Luong, Cho-Jui Hsieh, Yifeng Lu, et al. Symbolic discovery of optimization algorithms. *Advances in neural information processing systems*, 2024. (Cited on page 46)
- Dami Choi, Christopher J Shallue, Zachary Nado, Jaehoon Lee, Chris J Maddison, and George E Dahl. On empirical comparisons of optimizers for deep learning. *arXiv preprint arXiv:1910.05446*, 2019. (Cited on pages 1 and 2)
- Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, et al. Palm: Scaling language modeling with pathways. *Journal of Machine Learning Research*, 2023. (Cited on page 1)

- Enea Monzio Compagnoni, Tianlin Liu, Rustem Islamov, Frank Norbert Proske, Antonio Orvieto, and Aurelien Lucchi. Adaptive methods through the lens of SDEs: Theoretical insights on the role of noise. In *The Thirteenth International Conference on Learning Representations*, 2025. URL <https://openreview.net/forum?id=ww3CLRhF1v>. (Cited on page 3)
- Ashok Cutkosky and Harsh Mehta. Momentum improves normalized sgd. In *International conference on machine learning*. PMLR, 2020. (Cited on pages 2 and 3)
- Ashok Cutkosky and Francesco Orabona. Momentum-based variance reduction in non-convex sgd. *Advances in neural information processing systems*, 2019. (Cited on pages 2 and 3)
- Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*. Ieee, 2009. (Cited on page 9)
- Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale. In *International Conference on Learning Representations*, 2021. (Cited on pages 8 and 9)
- John Duchi, Elad Hazan, and Yoram Singer. Adaptive subgradient methods for online learning and stochastic optimization. *Journal of machine learning research*, 2011. (Cited on page 49)
- Ilyas Fatkhullin, Alexander Tyurin, and Peter Richtárik. Momentum provably improves error feedback! *Advances in Neural Information Processing Systems*, 2024. (Cited on page 2)
- Simon Foucart. Lecture 6: Matrix norms and spectral radii. *lecture notes for the course NSTP187 at Drexel University, Philadelphia, PA, Fall, 2012*, 2012. (Cited on page 35)
- Jingwen Fu, Bohan Wang, Huishuai Zhang, Zhizheng Zhang, Wei Chen, and Nanning Zheng. When and why momentum accelerates sgd: An empirical study. *arXiv preprint arXiv:2306.09000*, 2023. (Cited on page 2)
- Guillaume Garrigos and Robert M Gower. Handbook of convergence theorems for (stochastic) gradient methods. *arXiv preprint arXiv:2301.11235*, 2023. (Cited on pages 7, 8, 19, 21, 24, and 28)
- Euhanna Ghadimi, Hamid Reza Feyzmahdavian, and Mikael Johansson. Global convergence of the heavy-ball method for convex optimization. In *2015 European control conference (ECC)*, 2015. (Cited on page 7)
- Ian J. Goodfellow, Yoshua Bengio, and Aaron Courville. *Deep Learning*. MIT Press, 2016. (Cited on page 2)
- Robert Gower, Othmane Sebbouh, and Nicolas Loizou. Sgd for structured nonconvex functions: Learning rates, minibatching and interpolation. In *International Conference on Artificial Intelligence and Statistics*, 2021. (Cited on page 3)
- Robert M Gower, Guillaume Garrigos, Nicolas Loizou, Dimitris Oikonomou, Konstantin Mishchenko, and Fabian Schaipp. Analysis of an idealized stochastic polyak method and its application to black-box model distillation. *arXiv preprint arXiv:2504.01898*, 2025. (Cited on page 3)
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016. (Cited on page 8)
- Geoffrey Hinton, Nitish Srivastava, and Kevin Swersky. Neural networks for machine learning lecture 6a overview of mini-batch gradient descent. *Lecture notes*, 2012. (Cited on page 1)

- Tran Hoang, Qinzi Zhang, and Ashok Cutkosky. Empirical tests of optimization assumptions in deep learning. *arXiv preprint arXiv:2407.01825*, 2024. (Cited on page 6)
- Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. *Neural Computation*, 1997. (Cited on page 46)
- Jordan Hoffmann, Sebastian Borgeaud, Arthur Mensch, Elena Buchatskaya, Trevor Cai, Eliza Rutherford, Diego de Las Casas, Lisa Anne Hendricks, Johannes Welbl, Aidan Clark, Tom Hennigan, Eric Noland, Katie Millican, George van den Driessche, Bogdan Damoc, Aurelia Guy, Simon Osindero, Karen Simonyan, Erich Elsen, Jack W. Rae, Oriol Vinyals, and Laurent Sifre. Training compute-optimal large language models, 2022. URL <https://arxiv.org/abs/2203.15556>. (Cited on page 10)
- Rustem Islamov, Niccoló Ajroldi, Antonio Orvieto, and Aurelien Lucchi. Loss landscape characterization of neural networks without over-parametrization. In *Advances in Neural Information Processing Systems*, 2024a. (Cited on page 6)
- Rustem Islamov, Yuan Gao, and Sebastian U Stich. Near optimal decentralized optimization with compression and momentum tracking. *arXiv preprint arXiv:2405.20111*, 2024b. (Cited on page 2)
- Rustem Islamov, Samuel Horvath, Aurelien Lucchi, Peter Richtarik, and Eduard Gorbunov. Double momentum and error feedback for clipping with fast rates and differential privacy. *arXiv preprint arXiv: 2502.11682*, 2025. (Cited on page 2)
- Prateek Jain, Sham M Kakade, Rahul Kidambi, Praneeth Netrapalli, and Aaron Sidford. Accelerating stochastic gradient descent for least squares regression. In *Conference On Learning Theory*, 2018. (Cited on page 3)
- Samy Jelassi and Yuanzhi Li. Towards understanding how momentum improves generalization in deep learning. In *International Conference on Machine Learning*, 2022. (Cited on page 2)
- Xiaowen Jiang and Sebastian U Stich. Adaptive sgd with polyak stepsize and line-search: Robust convergence and variance reduction. *Advances in Neural Information Processing Systems*, 2024. (Cited on page 2)
- Keller Jordan, Yuchen Jin, Vlado Boza, Jiacheng You, Franz Cesista, Laker Newhouse, and Jeremy Bernstein. Muon: An optimizer for hidden layers in neural networks, 2024. URL <https://kellerjordan.github.io/posts/muon/>. (Cited on page 11)
- Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B. Brown, Benjamin Chess, Rewon Child, Scott Gray, Alec Radford, Jeffrey Wu, and Dario Amodei. Scaling laws for neural language models. *arXiv preprint arXiv: 2001.08361*, 2020. (Cited on page 45)
- Andrej Karpathy. char-rnn. <https://github.com/karpathy/char-rnn>, 2015. (Cited on page 46)
- Andrej Karpathy. Nanogpt. <https://github.com/karpathy/nanoGPT>, 2022. (Cited on pages 44 and 46)
- Diederik Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *International Conference on Learning Representations (ICLR)*, 2015. (Cited on pages 1 and 49)
- Alex Krizhevsky, Vinod Nair, and Geoffrey Hinton. Learning multiple layers of features from tiny images. *Scientific Report*, 2009. (Cited on page 8)
- Frederik Kunstner, Jacques Chen, Jonathan Wilder Lavington, and Mark Schmidt. Noise is not the main factor behind the gap between sgd and adam on transformers, but sign descent might be. In *The Eleventh International Conference on Learning Representations*, 2023. (Cited on page 3)

- Frederik Kunstner, Robin Yadav, Alan Milligan, Mark Schmidt, and Alberto Bietti. Heavy-tailed class imbalance and why adam outperforms gradient descent on language models. *arXiv preprint arXiv: 2402.19449*, 2024. (Cited on page 3)
- Kiwon Lee, Andrew Cheng, Elliot Paquette, and Courtney Paquette. Trajectory of mini-batch momentum: batch size saturation and convergence in high dimensions. *Advances in Neural Information Processing Systems*, 2022. (Cited on page 3)
- Haochuan Li, Alexander Rakhlin, and Ali Jadbabaie. Convergence of adam under relaxed assumptions. *Advances in Neural Information Processing Systems*, 2024. (Cited on page 1)
- Shen Li, Yanli Zhao, Rohan Varma, Omkar Salpekar, Pieter Noordhuis, Teng Li, Adam Paszke, Jeff Smith, Brian Vaughan, Pritam Damania, et al. Pytorch distributed: Experiences on accelerating data parallel training. *arXiv preprint arXiv:2006.15704*, 2020. (Cited on page 38)
- Shuang Li, William J Swartworth, Martin Takáč, Deanna Needell, and Robert M Gower. Sp2: A second order stochastic polyak method. *arXiv preprint arXiv:2207.08171*, 2022. (Cited on page 3)
- Yanli Liu, Yuan Gao, and Wotao Yin. An improved analysis of stochastic gradient descent with momentum. *Advances in Neural Information Processing Systems*, 2020. (Cited on page 7)
- Nicolas Loizou, Sharan Vaswani, Issam Hadj Laradji, and Simon Lacoste-Julien. Stochastic polyak step-size for sgd: An adaptive learning rate for fast convergence. In *International Conference on Artificial Intelligence and Statistics*, 2021. (Cited on pages 2, 3, 6, and 7)
- Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. *arXiv preprint arXiv: 1711.05101*, 2019. (Cited on page 39)
- Liangchen Luo, Yuanhao Xiong, Yan Liu, and Xu Sun. Adaptive gradient methods with dynamic bound of learning rate. *arXiv preprint arXiv:1902.09843*, 2019. (Cited on page 46)
- Jerry Ma and Denis Yarats. Quasi-hyperbolic momentum and adam for deep learning. *arXiv preprint arXiv:1810.06801*, 2018. (Cited on pages 2 and 3)
- Stephen Merity, Caiming Xiong, James Bradbury, and Richard Socher. Pointer sentinel mixture models. In *Proceedings of 4th International Conference on Learning Representations (ICLR 2016)*, 2016. (Cited on page 46)
- Tomas Mikolov, Martin Karafiát, Lukas Burget, Jan Cernocký, and Sanjeev Khudanpur. Recurrent neural network based language model. *Proceedings of the 11th Annual Conference of the International Speech Communication Association, INTERSPEECH 2010*, 2010. (Cited on page 46)
- Igor Molybog, Peter Albert, Moya Chen, Zachary DeVito, David Esiobu, Naman Goyal, Punit Singh Koura, Sharan Narang, Andrew Poulton, Ruan Silva, et al. A theory on adam instability in large-scale machine learning. *arXiv preprint arXiv:2304.09871*, 2023. (Cited on page 1)
- Lorenzo Noci, Sotiris Anagnostidis, Luca Biggio, Antonio Orvieto, Sidak Pal Singh, and Aurelien Lucchi. Signal propagation in transformers: Theoretical perspectives and the role of rank collapse. *Advances in Neural Information Processing Systems*, 2022. (Cited on page 3)
- Dimitris Oikonomou and Nicolas Loizou. Stochastic polyak step-sizes and momentum: Convergence guarantees and practical performance. *arXiv preprint arXiv:2406.04142*, 2024. (Cited on pages 3, 4, 6, 7, and 19)
- Sharir Or, Barak Peleg, and Yoav Shoham. The cost of training nlp models: A concise overview. *arXiv preprint arXiv:2004.08900*, 2020. (Cited on page 1)
- Antonio Orvieto and Lin Xiao. An adaptive stochastic gradient method with non-negative gauss-newton stepsizes. *arXiv preprint arXiv: 2407.04358*, 2024. (Cited on pages 1, 2, 3, 4, 6, and 19)

- Antonio Orvieto, Simon Lacoste-Julien, and Nicolas Loizou. Dynamics of sgd with stochastic polyak stepsizes: Truly adaptive variants and convergence to exact solution. *Advances in Neural Information Processing Systems*, 2022. (Cited on pages 2 and 6)
- Bo Pang and Lillian Lee. Seeing stars: Exploiting class relationships for sentiment categorization with respect to rating scales. In *Proceedings of the ACL*, 2005. (Cited on page 46)
- Adam Paszke, Sam Gross, Soumith Chintala, Gregory Chanan, Edward Yang, Zachary DeVito, Zeming Lin, Alban Desmaison, Luca Antiga, and Adam Lerer. Automatic differentiation in pytorch. *NIPS 2017 Workshop Autodiff*, 2017. (Cited on page 44)
- Guilherme Penedo, Hynek Kydlíček, Anton Lozhkov, Margaret Mitchell, Colin A Raffel, Leandro Von Werra, Thomas Wolf, et al. The fineweb datasets: Decanting the web for the finest text data at scale. *Advances in Neural Information Processing Systems*, 2024. (Cited on page 50)
- Boris T. Polyak. Some methods of speeding up the convergence of iteration methods. *Ussr computational mathematics and mathematical physics*, 1964. (Cited on pages 1, 3, and 34)
- Boris T Polyak. Introduction to optimization. *New York, Optimization Software*, 1987. (Cited on page 2)
- Alec Radford, Karthik Narasimhan, Tim Salimans, and Ilya Sutskever. Improving language understanding by generative pre-training. Technical report, OpenAI, 2019. (Cited on page 10)
- Samyam Rajbhandari, Jeff Rasley, Olatunji Ruwase, and Yuxiong He. Zero: Memory optimizations toward training trillion parameter models. In *SC20: International Conference for High Performance Computing, Networking, Storage and Analysis*, pages 1–16. IEEE, 2020. (Cited on pages 11 and 38)
- Peter Richtárik, Simone Maria Giancola, Dymitr Lubczyk, and Robin Yadav. Local curvature descent: Squeezing more curvature out of standard and polyak gradient descent. *arXiv preprint arXiv:2405.16574*, 2024. (Cited on page 3)
- Samer Saab, Shashi Phoha, Minghui Zhu, and Asok Ray. An adaptive polyak heavy-ball method. *Machine Learning*, 2022. (Cited on page 3)
- Mher Safaryan and Peter Richtárik. Stochastic sign descent methods: New algorithms and better theory. In *International Conference on Machine Learning*, 2021. (Cited on page 19)
- Fabian Schaipp, Ruben Ohana, Michael Eickenberg, Aaron Defazio, and Robert M. Gower. MoMo: Momentum models for adaptive learning rates. In *Proceedings of the 41st International Conference on Machine Learning*, 2024. (Cited on pages 1, 3, 4, 6, 9, and 44)
- Fabian Schaipp, Alexander Hägele, Adrien Taylor, Umut Simsekli, and Francis Bach. The surprising agreement between convex optimization theory and learning-rate scheduling for large model training. *arXiv preprint arXiv:2501.18965*, 2025. (Cited on page 6)
- Othmane Sebbouh, Robert M Gower, and Aaron Defazio. Almost sure convergence rates for stochastic gradient descent and stochastic heavy ball. In *Conference on Learning Theory*, 2021. (Cited on pages 7 and 19)
- Noam Shazeer. Glu variants improve transformer. *arXiv preprint arXiv: 2002.05202*, 2020. (Cited on page 44)
- Mohammad Shoeybi, Mostofa Patwary, Raul Puri, Patrick LeGresley, Jared Casper, and Bryan Catanzaro. Megatron-lm: Training multi-billion parameter language models using model parallelism. *arXiv preprint arXiv:1909.08053*, 2019. (Cited on pages 11 and 38)
- Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*, 2014. (Cited on page 45)

- J. W. Smith, J. E. Everhart, W. C. Dickson, W. C. Knowler, and R. S. Johannes. Using the adap learning algorithm to forecast the onset of diabetes mellitus. In *Symposium on Computer Applications and Medical Care*, 1988. (Cited on page 4)
- Daria Soboleva, Faisal Al-Khateeb, Robert Myers, Jacob R Steeves, Joel Hestness, and Nolan Dey. SlimPajama: A 627B token cleaned and deduplicated version of RedPajama, 2023. (Cited on pages 10 and 45)
- Jianlin Su, Yu Lu, Shengfeng Pan, Ahmed Murtadha, Bo Wen, and Yunfeng Liu. Roformer: Enhanced transformer with rotary position embedding, 2023. (Cited on page 44)
- Hugo Touvron, Matthieu Cord, Matthijs Douze, Francisco Massa, Alexandre Sablayrolles, and Hervé Jégou. Training data-efficient image transformers & distillation through attention. In *International conference on machine learning*, 2021. (Cited on page 1)
- Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurelien Rodriguez, Armand Joulin, Edouard Grave, and Guillaume Lample. Llama: Open and efficient foundation language models. *arXiv preprint arXiv: 2302.13971*, 2023. (Cited on pages 1 and 45)
- Bohan Wang, Jingwen Fu, Huishuai Zhang, Nanning Zheng, and Wei Chen. Closing the gap between the upper bound and lower bound of adam’s iteration complexity. *Advances in Neural Information Processing Systems*, 2024. (Cited on page 1)
- Jun-Kun Wang, Chi-Heng Lin, and Jacob D Abernethy. A modular analysis of provable acceleration via polyak’s momentum: Training a wide relu network and a deep linear network. In *International Conference on Machine Learning*, pages 10816–10827. PMLR, 2021. (Cited on page 35)
- Jun-Kun Wang, Chi-Heng Lin, Andre Wibisono, and Bin Hu. Provable acceleration of heavy ball beyond quadratics for a class of polyak-lojasiewicz functions when the non-convexity is averaged-out. In *International conference on machine learning*, 2022. (Cited on pages 3 and 7)
- Xiaoyu Wang, Mikael Johansson, and Tong Zhang. Generalized polyak step size for first order optimization with momentum. In *International Conference on Machine Learning*, 2023. (Cited on pages 3, 4, and 7)
- Rachel Ward, Xiaoxia Wu, and Leon Bottou. Adagrad stepsizes: Sharp convergence over nonconvex landscapes. *Journal of Machine Learning Research*, 2020. (Cited on page 20)
- Ross Wightman. Pytorch image models. <https://github.com/rwightman/pytorch-image-models>, 2019. (Cited on page 44)
- Ashia C Wilson, Rebecca Roelofs, Mitchell Stern, Nati Srebro, and Benjamin Recht. The marginal value of adaptive gradient methods in machine learning. *Advances in neural information processing systems*, 2017. (Cited on page 1)
- Mitchell Wortsman, Peter J Liu, Lechao Xiao, Katie Everett, Alex Alemi, Ben Adlam, John D Co-Reyes, Izzeddin Gur, Abhishek Kumar, Roman Novak, et al. Small-scale proxies for large-scale transformer training instabilities. *arXiv preprint arXiv:2309.14322*, 2023. (Cited on page 1)
- Lechao Xiao. Rethinking conventional wisdom in machine learning: From generalization to scaling. *arXiv preprint arXiv: 2409.15156*, 2024. (Cited on pages 38 and 41)
- Greg Yang, Edward J Hu, Igor Babuschkin, Szymon Sidor, Xiaodong Liu, David Farhi, Nick Ryder, Jakub Pachocki, Weizhu Chen, and Jianfeng Gao. Tensor programs v: Tuning large neural networks via zero-shot hyperparameter transfer. *arXiv preprint arXiv:2203.03466*, 2022. (Cited on page 2)

- Greg Yang, James B Simon, and Jeremy Bernstein. A spectral condition for feature learning. *arXiv preprint arXiv:2310.17813*, 2023. (Cited on page 2)
- Biao Zhang and Rico Sennrich. Root mean square layer normalization, 2019. (Cited on page 44)
- Guodong Zhang, Chaoqi Wang, Bowen Xu, and Roger Grosse. Three mechanisms of weight decay regularization. *arXiv preprint arXiv:1810.12281*, 2018. (Cited on page 39)
- Jingzhao Zhang, Tianxing He, Suvrit Sra, and Ali Jadbabaie. Why gradient clipping accelerates training: A theoretical justification for adaptivity. *arXiv preprint arXiv:1905.11881*, 2019. (Cited on page 3)
- Jingzhao Zhang, Sai Praneeth Karimireddy, Andreas Veit, Seungyeon Kim, Sashank Reddi, Sanjiv Kumar, and Suvrit Sra. Why are adaptive methods good for attention models? *Advances in Neural Information Processing Systems*, 2020. (Cited on page 3)
- Qi Zhang, Yi Zhou, and Shaofeng Zou. Convergence guarantees for rmsprop and adam in generalized-smooth non-convex optimization with affine noise variance. *arXiv preprint arXiv:2404.01436*, 2024a. (Cited on page 1)
- Susan Zhang, Stephen Roller, Naman Goyal, Mikel Artetxe, Moya Chen, Shuohui Chen, Christopher Dewan, Mona Diab, Xian Li, Xi Victoria Lin, et al. Opt: Open pre-trained transformer language models. *arXiv preprint arXiv:2205.01068*, 2022. (Cited on page 1)
- Yushun Zhang, Congliang Chen, Tian Ding, Ziniu Li, Ruoyu Sun, and Zhi-Quan Luo. Why transformers need adam: A hessian perspective. *arXiv preprint arXiv:2402.16788*, 2024b. (Cited on page 3)
- Juntang Zhuang, Tommy Tang, Yifan Ding, Sekhar C Tatikonda, Nicha Dvornek, Xenophon Papademetris, and James Duncan. Adabelief optimizer: Adapting stepsizes by the belief in observed gradients. In *Advances in Neural Information Processing Systems*, 2020. (Cited on page 46)

Appendix

Contents

A	Equivalent Formulations of NGN-M	18
B	Technical Lemmas and Definitions	19
C	Convergence of NGN-D	20
C.1	Convergence in General Non-convex Setting	22
C.2	Convergence under PL-condition	23
D	Convergence of NGN-M	24
E	Convergence of NGN-M with Decaying Step-size	28
F	Stability of NGN-M on a Simple Problem	33
G	How to Derive Diagonal NGN-based Step-size?	37
G.1	Design Comparison of NGN-MDv1 and NGN-MDv2	37
G.2	Computation Cost of NGN-MD	38
H	How to add weight decay to NGN-MDv1?	38
H.1	Combining NGN-MDv1 and Weight Decay Regularization	39
H.2	Empirical Validation of the Proposed Combinations	41
I	Additional Experiments on Toy Problems	42
I.1	Additional Experiments on the Problem with Many Minima	42
I.2	Comparison on Rosenbrock Function	42
I.3	Comparison on Quadratic Function with Theoretical Step-size	43
J	Additional Experiments and Training Details	44
J.1	Training Details	44
J.2	Comparison Algorithms that Support Momentum	45
J.3	Comparison of Algorithms that Support Momentum and Diagonal Step-size	46
J.4	Additional ImageNet Experiments	46
J.5	Additional Comparison against Lion, Adabelief, Adabound	46
J.6	Comparison of Adaptive Step-sizes of Adam, Momo-Adam, and NGN-MDv1	47
J.7	Extended Comparison of Momentum-based Algorithms on NLP Tasks	48
J.8	Comparison of Algorithms with Diagonal Step-size	49
J.9	Effective Step-size of NGN-M, Momo, NGN-MDv1, and Momo-Adam	49
J.10	Effective Updates in Training Language Models	50
J.11	Training Dynamics in Training Language Models	50
J.12	Ablation Study of Momentum Parameters	50

A Equivalent Formulations of NGN-M

We remind that the iterates of NGN-M are the following

$$\begin{aligned}x^{k+1} &= x^k - (1 - \beta)\gamma_k \nabla f_{S_k}(x^k) + \beta(x^k - x^{k-1}) \\ &= x^k - (1 - \beta) \frac{c}{1 + \frac{c}{2f_{S_k}(x^k)} \|\nabla f_{S_k}(x^k)\|^2} \nabla f_{S_k}(x^k) + \beta(x^k - x^{k-1}).\end{aligned}$$

We can rewrite the update rule using Iterative-Moving Average (IMA) approach presented in Proposition 1.6, [Sebbouh et al. \[2021\]](#).

Lemma 1 (Proposition C.8 [[Oikonomou and Loizou, 2024](#)], Lemma 7.3 in [[Garrigos and Gower, 2023](#)]). The iterates $\{x^k\}$ generated by NGN-M are equivalent to the sequence $\{x^k\}$ generated by IMA update

$$z^{k+1} = z^k - \gamma_k \nabla f_{S_k}(x^k), \quad x^{k+1} = \frac{\lambda}{1+\lambda} x^k + \frac{1}{1+\lambda} z^{k+1}, \quad (5)$$

where

$$\beta = \frac{\lambda}{1+\lambda}, \quad z^{k+1} = x^{k+1} + \lambda(x^{k+1} - x^k), \quad \text{and} \quad x^{-1} = z^0 = x^0. \quad (6)$$

Proof. Let the sequences $\{x^k\}$ and $\{z^k\}$ be defined according to Equation (5). Let β be defined as $\frac{\lambda}{1+\lambda}$. Then we have

$$\begin{aligned} x^{k+1} &= \frac{\lambda}{1+\lambda} x^k + \frac{1}{1+\lambda} z^{k+1} \\ &= \frac{\lambda}{1+\lambda} x^k + \frac{1}{1+\lambda} (z^k - \gamma_k \nabla f_{S_k}(x^k)) \\ &= \frac{\lambda}{1+\lambda} x^k + \frac{1}{1+\lambda} ((1+\lambda)x^k - \lambda x^{k-1} - \gamma_k \nabla f_{S_k}(x^k)) \\ &= x^k - \frac{1}{1+\lambda} \gamma_k \nabla f_{S_k}(x^k) + \frac{\lambda}{1+\lambda} (x^k - x^{k-1}). \end{aligned}$$

It remains to use equation 6 as we have $\beta = \frac{\lambda}{1+\lambda}$ and $1 - \beta = 1 - \frac{\lambda}{1+\lambda} = \frac{1}{1+\lambda}$. □

B Technical Lemmas and Definitions

Definition 1. We say that the function ϕ admits \mathbf{L} -smooth with parameters $\mathbf{L} := (L_1, \dots, L_d)$, $L_j \geq 0 \forall j \in [d]$, if the following inequality holds for all $x, h \in \mathbb{R}^d$

$$\phi(x+h) \leq \phi(x) + \langle \nabla \phi(x), h \rangle + \frac{1}{2} h^\top \mathbf{L} h. \quad (7)$$

Remark 1. If we set for all $j \in [d]$ $L_j := L$ then Definition 1 reduces to standard L -smoothness.

This assumption is typically used in the context of coordinate adaptive algorithms such as SignSGD [[Bernstein et al., 2018](#), [Safaryan and Richtárik, 2021](#)].

Definition 2. The function $\phi: \mathbb{R}^d \rightarrow \mathbb{R}$ satisfies PL -condition with constant $\mu > 0$ if for all $x, y \in \mathbb{R}^d$ we have

$$\|\nabla f(x)\|^2 \geq 2\mu(f(x) - f^*). \quad (8)$$

Assumption 3. We assume that the coordinate-wise variance of the stochastic estimator is bounded, i.e. for all $x \in \mathbb{R}^d$ and $j \in [d]$ we have

$$\mathbb{E}_S [|(\nabla_j f_S(x) - \nabla_j f(x))|^2] \leq \sigma_j^2. \quad (9)$$

Lemma 2 (Lemma 4.9 from [[Orvieto and Xiao, 2024](#)]). Let each f_i be L -smooth for all i , then the step-size of NGN satisfies

$$\gamma_k \in \left[\frac{c}{1+cL}, c \right]. \quad (10)$$

Lemma 3 (Lemma 4.2 from [[Orvieto and Xiao, 2024](#)]). Let each f_i be L -smooth for all i , then the iterates of NGN satisfy

$$\gamma_k^2 \|\nabla f_{S_k}(x^k)\|^2 \leq \frac{4cL}{1+2cL} \gamma_k (f_{S_k}(x^k) - f_{S_k}^*) + \frac{2c^2L}{1+cL} \max \left\{ \frac{2cL-1}{2cL+1}, 0 \right\} f_{S_k}^*. \quad (11)$$

Lemma 4 (Gradient Upper Bound). Let $\phi: \mathbb{R}^d \rightarrow \mathbb{R}$ satisfy Definition 1. Then, for all $x \in \mathbb{R}^d$ and all $j \in [d]$ we have

$$2L_j(f(x) - f^*) \geq (\nabla_j f(x))^2. \quad (12)$$

Proof. From Definition 1 we have

$$f^* = \min_{x \in \mathbb{R}^d} f(x) \leq \min_{h_j \in \mathbb{R}} f(x + h_j e_j) \leq f(x) + \min_{h_j \in \mathbb{R}} \left[\nabla_j f(x) h_j + \frac{L_j}{2} h_j^2 \right].$$

Now we can explicitly compute the minimum in the right-hand side. The optimal value is achieved at

$$h_j^* := -\frac{1}{L_j} \nabla_j f(x),$$

therefore,

$$\begin{aligned} f^* &\leq f(x) + \nabla_j f(x) h_j^* + \frac{L_j}{2} (h_j^*)^2 \\ &= f(x) - \frac{1}{L_j} (\nabla_j f(x))^2 + \frac{1}{2L_j} (\nabla_j f(x))^2 \\ &= f(x) - \frac{1}{2L_j} (\nabla_j f(x))^2, \end{aligned}$$

which equivalent to the statement of the lemma. \square

C Convergence of NGN-D

First, we provide NGN-D pseudocode and the main convergence results.

Algorithm 3 NGN-D

- 1: **Input:** $x^0 \in \mathbb{R}^d$, step-size parameter $c > 0$
 - 2: **for** $k = 0, 1, \dots, K - 1$ **do**
 - 3: Sample a batch $S_k \subseteq [n]$ and compute f_{S_k} and $\nabla f_{S_k}(x^k)$
 - 4: Compute $\gamma_k^{(j)} = \frac{c}{1 + \frac{c}{2f_{S_k}(x^k)} (\nabla_j f_{S_k}(x^k))^2}$
 - 5: Update

$$x_{(j)}^{k+1} = x_{(j)}^k - \gamma_k^{(j)} \nabla_j f_{S_k}(x^k)$$
 - 6: **end for**
-

Theorem 2. Let each f_i satisfies Definition 1. Assume that Assumption 3 holds. Then the iterates of NGN-D (Alg. 3) with step-size parameters $\{c_j\}_{j=1}^d$ such that $c_j \leq 1/2L_j$ satisfy

$$\min_{0 \leq k < K} \mathbb{E} [\|\nabla f(x^k)\|^2] \leq \frac{12(f(x^0) - f^*)}{c_{\min} K} + \frac{1}{c_{\min}} \sum_{j=1}^d 18L_j c_j^2 \sigma_j^2, \quad (13)$$

where $c_{\min} := \min_{j \in [d]} c_j$. Moreover, if $c_j = \mathcal{O}(\varepsilon^2)$ for all $j \in [d]$ then after $K = \mathcal{O}(\varepsilon^{-4})$ we obtain $\min_{0 \leq k < K} \mathbb{E} [\|\nabla f(x^k)\|^2] \leq \mathcal{O}(\varepsilon^2)$.

NGN-D converges with classic rate $\mathcal{O}(1/\sqrt{K})$ similar to Adagrad [Ward et al., 2020]. We highlight that, to the best of our knowledge, NGN-D is the first algorithm that uses diagonal Polyak-type stepsize and converges under standard smoothness and bounded variance assumptions without requirements of bounded gradients and interpolation.

Theorem 3. Let f satisfies PL-condition and each f_i satisfies Definition 1. Assume that Assumption 3 holds. Then the iterates of NGN-D (Alg. 3) with step-size parameters $\{c_j\}_{j=1}^d$ such that $c_j \leq \min\{1/2L_j, 6/\mu\}$ satisfy

$$\mathbb{E} [f(x^K) - f^*] \leq (1 - \mu c_{\min}/6)^K (f(x^0) - f^*) + \frac{9}{\mu c_{\min}} \sum_{j=1}^d L_j c_j^2 \sigma_j^2, \quad (14)$$

where $c_{\min} := \min_{j \in [d]} c_j$. Moreover, if $c_j = \mathcal{O}(\varepsilon)$ for all $j \in [d]$ then after $K = \max\{\mathcal{O}(\varepsilon^{-1}), \mathcal{O}(1)\} \log \varepsilon^{-1}$ iterations we obtain $\mathbb{E} [f(x^K) - f^*] \leq \mathcal{O}(\varepsilon)$.

To the best of our knowledge, this is the first result of the convergence of the Polyak-like step-size algorithm under the PL-condition. The convergence guarantees are similar to that of SGD [Garrigos and Gower, 2023].

Now we are ready to derive the step-size bounds.

Lemma 5 (Step-size Bounds). Let $f_{S_k}(x): \mathbb{R}^d \rightarrow \mathbb{R}$ be a stochastic loss of batch S_k at iteration k . Let $f_{S_k}(x)$ satisfy Definition equation 1. Consider γ_j^k as in NGN-D (Algorithm 3), then we have

$$\gamma_j^k \in \left[\frac{c_j}{1 + c_j L_j}, c_j \right]. \quad (15)$$

Proof. From Lemma 4 we have $2L_j(f_{S_k}(x^k) - f_{S_k}^*) \geq (\nabla_j f_{S_k}(x^k))^2$. Since we assume that each $f_{S_k}^* \geq 0$, then $2L_j f_{S_k}(x^k) \geq (\nabla_j f_{S_k}(x^k))^2$, or equivalently,

$$0 \leq \frac{(\nabla_j f_{S_k}(x))^2}{2f_{S_k}(x)} \leq L_j.$$

Therefore, for all $j \in [d]$ we have

$$\gamma_j^k = \frac{c_j}{1 + \frac{c_j}{2f_{S_k}(x^k)} (\nabla_j f_{S_k}(x^k))^2} \leq \frac{c_j}{1} = c_j,$$

and

$$\gamma_j^k = \frac{c_j}{1 + \frac{c_j}{2f_{S_k}(x^k)} (\nabla_j f_{S_k}(x^k))^2} \geq \frac{c_j}{1 + c_j L_j},$$

that concludes the proof. \square

Lemma 6 (Fundamental Equality). Consider γ_j^k as in NGN-D (Algorithm 3). Then the following equality holds

$$\gamma_j^k (\nabla_j f_{S_k}(x^k))^2 = 2 \left(\frac{c_j - \gamma_j^k}{c_j} \right) f_{S_k}(x^k). \quad (16)$$

Proof. From NGN-D (Algorithm 3) we have

$$\left(1 + \frac{c_j}{2f_{S_k}(x^k)} (\nabla_j f_{S_k}(x^k))^2 \right) \gamma_j^k = c_j,$$

which one can rewrite as

$$\frac{c_j}{2f_{S_k}(x^k)} (\nabla_j f_{S_k}(x^k))^2 \gamma_j^k = c_j - \gamma_j^k.$$

It is left to divide both sides by $\frac{2f_{S_k}(x^k)}{c_j}$. \square

C.1 Convergence in General Non-convex Setting

Theorem 2. Let each f_i satisfies Definition 1. Assume that Assumption 3 holds. Then the iterates of NGN-D (Alg. 3) with step-size parameters $\{c_j\}_{j=1}^d$ such that $c_j \leq 1/2L_j$ satisfy

$$\min_{0 \leq k < K} \mathbb{E} [\|\nabla f(x^k)\|^2] \leq \frac{12(f(x^0) - f^*)}{c_{\min} K} + \frac{1}{c_{\min}} \sum_{j=1}^d 18L_j c_j^2 \sigma_j^2, \quad (13)$$

where $c_{\min} := \min_{j \in [d]} c_j$. Moreover, if $c_j = \mathcal{O}(\varepsilon^2)$ for all $j \in [d]$ then after $K = \mathcal{O}(\varepsilon^{-4})$ we obtain $\min_{0 \leq k < K} \mathbb{E} [\|\nabla f(x^k)\|^2] \leq \mathcal{O}(\varepsilon^2)$.

Proof. First, we write separable Definition 1

$$\begin{aligned} f(x^{k+1}) - f(x^k) &= f\left(x^k - \sum_{j=1}^d \gamma_j^k \nabla_j f_{S_k}(x^k) e_j\right) - f(x^k) \\ &\leq -\sum_{j=1}^d \nabla_j f(x^k) \cdot \gamma_j^k \nabla_j f_{S_k}(x^k) + \frac{1}{2} \sum_{j=1}^d L_j (\gamma_j^k \nabla_j f_{S_k}(x^k))^2 \\ &\leq -\sum_{j=1}^d \nabla_j f(x^k) \cdot \gamma_j^k \nabla_j f_{S_k}(x^k) + \frac{1}{2} \sum_{j=1}^d L_j \sigma_j^2 (\nabla_j f_{S_k}(x^k))^2. \end{aligned} \quad (17)$$

Note that both γ_j^k and $\nabla_j f_{S_k}(x^k)$ depend on the realization S_k , thus we can not directly apply conditional expectation with respect to x^k , as in this case we would have to analyze the product $\gamma_j^k \nabla_j f_{S_k}(x^k)$. Given bounds of the step-size γ_j^k from Lemma 5, we can write the step-size as follows

$$\gamma_j^k = \frac{c_j}{1 + c_j L_j} + \nu_j^k \frac{c_j^2 L_j}{1 + c_j L_j},$$

where $\nu_j^k \in [0, 1]$ is a random variable. Varying the value of ν_j^k from 0 to 1 we cover the whole range of γ_j^k . Thus, we continue as follows

$$\begin{aligned} &-\gamma_j^k \nabla_j f(x^k) \nabla_j f_{S_k}(x^k) \\ &= -\frac{c_j}{1 + c_j L_j} \nabla_j f(x^k) \nabla_j f_{S_k}(x^k) - \frac{c_j^2 L_j}{1 + c_j L_j} \nu_j^k \nabla_j f(x^k) \nabla_j f_{S_k}(x^k) \\ &\leq -\frac{c_j}{1 + c_j L_j} \nabla_j f(x^k) \nabla_j f_{S_k}(x^k) + \frac{c_j^2 L_j}{1 + c_j L_j} |\nu_j^k| \cdot |\nabla_j f(x^k) \nabla_j f_{S_k}(x^k)| \\ &\leq -\frac{c_j}{1 + c_j L_j} \nabla_j f(x^k) \nabla_j f_{S_k}(x^k) + \frac{c_j^2 L_j}{1 + c_j L_j} \cdot |\nabla_j f(x^k) \nabla_j f_{S_k}(x^k)|. \end{aligned}$$

Now we use the inequality $|ab| \leq \frac{1}{2}a^2 + \frac{1}{2}b^2 + \frac{1}{2}|a - b|^2$, and derive

$$\begin{aligned} 2\mathbb{E}_k [|\nabla_j f(x^k) \nabla_j f_{S_k}(x^k)|] &\leq |\nabla_j f(x^k)|^2 + \mathbb{E}_k [|\nabla_j f_{S_k}(x^k)|^2] + \mathbb{E}_k [|\nabla_j f(x^k) - \nabla_j f_{S_k}(x^k)|^2] \\ &\leq 2|\nabla_j f(x^k)|^2 + 2\mathbb{E}_k [|\nabla_j f(x^k) - \nabla_j f_{S_k}(x^k)|^2] \\ &\leq 2|\nabla_j f(x^k)|^2 + 2\sigma_j^2. \end{aligned}$$

Therefore, we get

$$\begin{aligned} -\mathbb{E}_k [\gamma_j^k \nabla_j f(x^k) \nabla_j f_{S_k}(x^k)] &\leq -\frac{c_j}{1 + c_j L_j} |\nabla_j f(x^k)|^2 + \frac{c_j^2 L_j}{1 + c_j L_j} (|\nabla_j f(x^k)|^2 + \sigma_j^2) \\ &= -c_j \left(\frac{1 - c_j L_j}{1 + c_j L_j} \right) |\nabla_j f(x^k)|^2 + \frac{c_j^2 L_j}{1 + c_j L_j} \sigma_j^2. \end{aligned} \quad (18)$$

We plug in equation 18 into equation 17 and get

$$\begin{aligned}\mathbb{E}_k [f(x^{k+1})] - f(x^k) &\leq -\sum_{j=1}^d \left(\mathbb{E}_k [\gamma_j^k \nabla_j f(x^k) \nabla_j f_{S_k}(x^k)] + \frac{L_j c_j^2}{2} \mathbb{E}_k [|\nabla_j f_{S_k}(x^k)|^2] \right) \\ &\leq \sum_{j=1}^d \left(\left[-c_j \left(\frac{1 - c_j L_j}{1 + c_j L_j} \right) + \frac{L_j c_j^2}{2} \right] |\nabla_j f(x^k)|^2 \right. \\ &\quad \left. + \left[\frac{c_j^2 L_j}{1 + c_j L_j} + \frac{L_j c_j^2}{2} \right] \sigma_j^2 \right).\end{aligned}$$

If $c_j \leq \frac{1}{2L_j}$, we get

$$\mathbb{E}_k [f(x^{k+1})] - f(x^k) \leq \sum_{j=1}^d \left(-\frac{c_j}{12} |\nabla_j f(x^k)|^2 + \frac{3L_j c_j^2}{2} \sigma_j^2 \right).$$

□

We continue as follows

$$\mathbb{E}_k [f(x^{k+1})] - f(x^k) \leq -\frac{c_{\min}}{12} \|\nabla f(x^k)\|^2 + \sum_{j=1}^d \frac{3L_j c_j^2}{2} \sigma_j^2. \quad (19)$$

Taking full expectation and unrolling the recursion above for all iterations $\{0, \dots, K-1\}$. Thus, we obtain

$$\min_{0 \leq k < K} \mathbb{E} [\|\nabla f(x^k)\|^2] \leq \frac{1}{K} \sum_{k=0}^{K-1} \mathbb{E} [\|\nabla f(x^k)\|^2] \leq \frac{12}{c_{\min} K} (f(x^0) - f^*) + \frac{18}{c_{\min}} \sum_{j=1}^d L_j c_j^2 \sigma_j^2.$$

If we choose each $c_j = \frac{c_{0,j}}{\sqrt{K}}$ such that $c_{0,j} \leq \frac{1}{2L_j}$ we ensure that $c_j \leq \frac{1}{2L_j}$ as well. Plugging this step-size into the bound we get

$$\begin{aligned}\min_{0 \leq k < K} \mathbb{E} [\|\nabla f(x^k)\|^2] &\leq \frac{12}{\frac{c_{0,\min}}{\sqrt{K}} K} (f(x^0) - f^*) + \frac{18}{\frac{c_{0,\min}}{\sqrt{K}}} \sum_{j=1}^d L_j \sigma_j^2 \frac{c_{0,j}^2}{K} \\ &\leq \frac{12}{c_{0,\min} \sqrt{K}} (f(x^0) - f^*) + \frac{18}{c_{0,\min} \sqrt{K}} \sum_{j=1}^d L_j \sigma_j^2 c_{0,j}^2,\end{aligned}$$

where $c_{0,\min} := \min_{j \in [d]} c_{0,j}$. If we choose $K = \mathcal{O}(\varepsilon^{-4})$ we get that

$$\min_{0 \leq k < K} \mathbb{E} [\|\nabla f(x^k)\|^2] = \mathcal{O}(1/\sqrt{K}) = \mathcal{O}(\varepsilon^2).$$

C.2 Convergence under PL-condition

Theorem 3. *Let f satisfies PL-condition and each f_i satisfies Definition 1. Assume that Assumption 3 holds. Then the iterates of NGN-D (Alg. 3) with step-size parameters $\{c_j\}_{j=1}^d$ such that $c_j \leq \min\{1/2L_j, 6/\mu\}$ satisfy*

$$\mathbb{E} [f(x^K) - f^*] \leq (1 - \mu c_{\min}/6)^K (f(x^0) - f^*) + \frac{9}{\mu c_{\min}} \sum_{j=1}^d L_j c_j^2 \sigma_j^2, \quad (14)$$

where $c_{\min} := \min_{j \in [d]} c_j$. Moreover, if $c_j = \mathcal{O}(\varepsilon)$ for all $j \in [d]$ then after $K = \max\{\mathcal{O}(\varepsilon^{-1}), \mathcal{O}(1)\} \log \varepsilon^{-1}$ iterations we obtain $\mathbb{E} [f(x^K) - f^*] \leq \mathcal{O}(\varepsilon)$.

Proof. We obtain equation 19 and use Definition 2

$$\begin{aligned}\mathbb{E}_k \left[f(x^{k+1}) \right] - f(x^k) &\leq -\frac{c_{\min}}{12} \|\nabla f(x^k)\|^2 + \sum_{j=1}^d \frac{3L_j c_j^2}{2} \sigma_j^2 \\ &\leq -\frac{\mu c_{\min}}{6} (f(x^k) - f^*) + \sum_{j=1}^d \frac{3L_j c_j^2}{2} \sigma_j^2\end{aligned}$$

Subtracting f^* from both sides of the inequality above and taking full expectation we obtain

$$\mathbb{E} \left[f(x^{k+1}) - f^* \right] \leq (1 - \mu c_{\min}/6) \mathbb{E} \left[f(x^k) - f^* \right] + \sum_{j=1}^d \frac{3L_j c_j^2}{2} \sigma_j^2.$$

Unrolling the recursion above for $\{0, \dots, K-1\}$ iterations we derive

$$\mathbb{E} \left[f(x^K) - f^* \right] \leq (1 - \mu c_{\min}/6)^K (f(x^0) - f^*) + \frac{1}{c_{\min}} \sum_{j=1}^d \underbrace{\frac{9L_j \sigma_j^2}{\mu}}_{A_j} c_j^2.$$

Now we follow the proof of Lemma A.3 in Garrigos and Gower [2023]. Let us choose $c_j = \min\{1/2L_j, \varepsilon/2dA_j\}$. Together with the choice of $K \geq \max_{j \in [d]} \max \left\{ \frac{1}{\varepsilon} \frac{12A_j}{\mu}, \frac{12L_j}{\mu} \right\} \log \frac{2(f(x^0) - f^*)}{\varepsilon}$ we get

$$(1 - \mu c_{\min}/6)^K (f(x^0) - f^*) \leq \frac{\varepsilon}{2}.$$

Now we have two cases:

1. c_{\min} does not depend on ε , then we have

$$\frac{1}{c_{\min}} A_j c_j^2 \leq \mathcal{O}(\varepsilon^2).$$

2. c_{\min} does depend on ε , i.e. $c_{\min} = \mathcal{O}(\varepsilon)$, then we have

$$\frac{1}{c_{\min}} A_j c_j^2 \leq \mathcal{O}(\varepsilon).$$

Therefore, combining all together we get

$$\mathbb{E} \left[f(x^K) - f^* \right] \leq \mathcal{O}(\varepsilon)$$

after $K \geq \max_{j \in [d]} \max \left\{ \frac{1}{\varepsilon} \frac{12A_j}{\mu}, \frac{12L_j}{\mu} \right\} \log \frac{2(f(x^0) - f^*)}{\varepsilon}$ iterations. □

D Convergence of NGN-M

Theorem 1. *Let Assumptions 1, 2 hold. Let the step-size hyperparameter $c > 0$ and the momentum parameter $\beta = \frac{\lambda}{1+\lambda}$ be constants where $\lambda \leq \min\{cL, 0.5(1+cL)^{-1}(1+2cL)^{-1}\}$. Then the iterates of NGN-M (Alg. 1) satisfy*

$$\mathbb{E} \left[f(\bar{x}^{K-1}) - f(x^*) \right] \leq \frac{\|x^0 - x^*\|^2 (1+2cL)^2}{cK} + 8cL(1+2cL)^2 \sigma_{\text{int}}^2 + 2cL \max\{2cL - 1, 0\} \sigma_{\text{pos}}^2,$$

where \bar{x}^{K-1} is chosen uniformly at random from $\{x^0, \dots, x^{K-1}\}$. Moreover, if we set $c = \mathcal{O}(1/\sqrt{K})$ then we obtain $\mathbb{E} \left[f(\bar{x}^{K-1}) - f(x^*) \right] \leq \mathcal{O}(1/\sqrt{K})$.

Remark 2. In fact, if $\lambda \leq \frac{1}{(1+cL)(1+2cL)}$, then it implies that $\lambda \leq \frac{1}{cL}$ because $\frac{1}{x} > \frac{1}{(1+x)(1+2x)}$ for any $x > 0$.

Proof. To prove the convergence of NGN-M we consider IMA formulation Equation (5):

$$x^{-1} = z^0 = x^0, \quad z^{k+1} = z^k - \gamma_k \nabla f_{S_k}(x^k), \quad x^{k+1} = \frac{\lambda}{1+\lambda} x^k + \frac{1}{1+\lambda} z^{k+1},$$

where $\beta = \frac{\lambda}{1+\lambda}$, $z^{k+1} = x^{k+1} + \lambda(x^{k+1} - x^k)$.

At iteration $k = 0$ we have

$$z^1 = z^0 - \gamma_0 \nabla f_{S_0}(x^0) = x^0 - \gamma_0 \nabla f_{S_0}(x^0).$$

Therefore, we get

$$\begin{aligned} \|z^1 - x^*\|^2 &= \|z^0 - x^*\|^2 - 2\gamma_0 \langle \nabla f_{S_0}(x^0), z^0 - x^* \rangle + \gamma_0^2 \|\nabla f_{S_0}(x^0)\|^2 \\ &\stackrel{\text{Lem. 3}}{\leq} \|z^0 - x^*\|^2 - 2\gamma_0 \langle \nabla f_{S_0}(x^0), x^0 - x^* \rangle + \frac{4cL}{1+2cL} \gamma_0 (f_{S_0}(x^0) - f_{S_0}^*) \\ &\quad + \frac{2c^2L}{1+cL} \max \left\{ \frac{2cL-1}{2cL+1}, 0 \right\} f_{S_0}^*. \end{aligned} \quad (20)$$

Let $\gamma_0 = \rho + \tilde{\gamma}_0$ where $\rho = \frac{c}{(1+cL)(1+2cL)}$. Then we have

$$\begin{aligned} \tilde{\gamma}_0 &= \gamma_0 - \rho \\ &\stackrel{\text{Lem. 2}}{\leq} c - \frac{c}{(1+cL)(1+2cL)} \\ &= c \frac{1+3cL+2c^2L^2-1}{(1+cL)(1+2cL)} \\ &= c^2L \frac{3+3cL}{(1+cL)(1+2cL)} \\ &= \frac{3c^2L}{1+2cL}. \end{aligned}$$

Using the above we continue from (20)

$$\begin{aligned} \|z^1 - x^*\|^2 &\stackrel{\text{conv.}}{\leq} \|z^0 - x^*\|^2 - 2\gamma_0 (f_{S_0}(x^0) - f_{S_0}(x^*)) + \frac{4cL}{1+2cL} \gamma_0 (f_{S_0}(x^0) - f_{S_0}^*) \\ &\quad + \frac{2c^2L}{1+cL} \max \left\{ \frac{2cL-1}{2cL+1}, 0 \right\} f_{S_0}^* \\ &\leq \|z^0 - x^*\|^2 - 2\rho (f_{S_0}(x^0) - f_{S_0}(x^*)) - 2\tilde{\gamma}_0 (f_{S_0}(x^0) - f_{S_0}^*) + 2\tilde{\gamma}_0 (f_{S_0}(x^*) - f_{S_0}^*) \\ &\quad + \frac{4cL}{1+2cL} \gamma_0 (f_{S_0}(x^0) - f_{S_0}^*) + \frac{2c^2L}{1+cL} \max \left\{ \frac{2cL-1}{2cL+1}, 0 \right\} f_{S_0}^* \\ &= \|z^0 - x^*\|^2 - 2\rho (f_{S_0}(x^0) - f_{S_0}(x^*)) - 2 \left(\gamma_0 - \rho - \frac{2cL}{1+2cL} \gamma_0 \right) (f_{S_0}(x^0) - f_{S_0}^*) \\ &\quad + 2\tilde{\gamma}_0 (f_{S_0}(x^*) - f_{S_0}^*) + \frac{2c^2L}{1+cL} \max \left\{ \frac{2cL-1}{2cL+1}, 0 \right\} f_{S_0}^*. \end{aligned} \quad (21)$$

Here we have

$$\begin{aligned} \gamma_0 - \rho - \frac{2cL}{1+2cL} \gamma_0 &= \frac{1}{1+2cL} \gamma_0 - \rho \\ &= \frac{1}{1+2cL} \gamma_0 - \frac{c}{(1+cL)(1+2cL)} \\ &\stackrel{\text{Lem. 2}}{\geq} \frac{1}{1+2cL} \frac{c}{1+cL} - \frac{c}{(1+cL)(1+2cL)} \\ &= 0, \end{aligned}$$

$\tilde{\gamma}_0 \leq \frac{3c^2L}{1+2cL}$, and $f_{S_0}(x^0) - f_{S_0}^* \geq 0$. Hence, we get

$$\begin{aligned} \|z^1 - x^*\|^2 &\leq \|z^0 - x^*\|^2 - 2\rho(f_{S_0}(x^0) - f_{S_0}^*) + \frac{6c^2L}{1+2cL}(f_{S_0}(x^*) - f_{S_0}^*) \\ &\quad + \frac{2c^2L}{1+cL} \max\left\{\frac{2cL-1}{2cL+1}, 0\right\} f_{S_0}^*. \end{aligned}$$

Rearranging terms and taking expectation we get

$$\begin{aligned} 2\rho\mathbb{E}[f(x^0) - f(x^*)] &\leq \mathbb{E}[\|z^1 - x^*\|^2] - \|z^0 - x^*\|^2 + \frac{6c^2L}{1+2cL}\sigma_{\text{int}}^2 \\ &\quad + \frac{2c^2L}{1+cL} \max\left\{\frac{2cL-1}{2cL+1}, 0\right\} \sigma_{\text{pos}}^2. \end{aligned} \quad (22)$$

Next, for $k > 0$ we can use the relation $z^k = x^k + \lambda(x^k - x^{k-1})$. We expand $\|z^{k+1} - x^*\|^2$

$$\begin{aligned} \|z^{k+1} - x^*\|^2 &= \|z^k - x^*\|^2 - 2\gamma_k \langle \nabla f_{S_k}(x^k), z^k - x^* \rangle + \gamma_k^2 \|\nabla f_{S_k}(x^k)\|^2 \\ &\stackrel{\text{Lem. 1}}{=} \|z^k - x^*\|^2 - 2\gamma_k \langle \nabla f_{S_k}(x^k), x^k - x^* \rangle - 2\gamma_k \lambda \langle \nabla f_{S_k}(x^k), x^k - x^{k-1} \rangle \\ &\quad + \gamma_k^2 \|\nabla f_{S_k}(x^k)\|^2 \\ &\stackrel{\text{conv.}}{\leq} \|z^k - x^*\|^2 - 2\gamma_k(f_{S_k}(x^k) - f_{S_k}^*) - 2\gamma_k \lambda(f_{S_k}(x^k) - f_{S_k}(x^{k-1})) \\ &\quad + \gamma_k^2 \|\nabla f_{S_k}(x^k)\|^2 \\ &\stackrel{\text{Lem. 3}}{\leq} \|z^k - x^*\|^2 - 2\gamma_k(f_{S_k}(x^k) - f_{S_k}^*) - 2\gamma_k \lambda(f_{S_k}(x^k) - f_{S_k}(x^{k-1})) \\ &\quad + \frac{4cL}{1+2cL} \gamma_k(f_{S_k}(x^k) - f_{S_k}^*) + \frac{2c^2L}{1+cL} \max\left\{\frac{2cL-1}{2cL+1}, 0\right\} f_{S_k}^*. \end{aligned}$$

Let $\gamma_k = \rho + \tilde{\gamma}_k$, where $\rho, \tilde{\gamma}_k \geq 0$, and ρ is a constant step-size independent of S_k which will be defined later. Therefore, we have

$$\begin{aligned} \|z^{k+1} - x^*\|^2 &\leq \|z^k - x^*\|^2 - 2\rho(f_{S_k}(x^k) - f_{S_k}^*) - 2\tilde{\gamma}_k(f_{S_k}(x^k) - f_{S_k}^*) \\ &\quad - 2\gamma_k \lambda_k(f_{S_k}(x^k) - f_{S_k}^*) + 2\gamma_k \lambda(f_{S_k}(x^{k-1}) - f_{S_k}^*) \\ &\quad + \frac{4cL}{1+2cL} \gamma_k(f_{S_k}(x^k) - f_{S_k}^*) + \frac{2c^2L}{1+cL} \max\left\{\frac{2cL-1}{2cL+1}, 0\right\} f_{S_k}^* \\ &= \|z^k - x^*\|^2 - 2\rho(f_{S_k}(x^k) - f_{S_k}^*) - 2\tilde{\gamma}_k(f_{S_k}(x^k) - f_{S_k}^*) + 2\tilde{\gamma}_k(f_{S_k}(x^*) - f_{S_k}^*) \\ &\quad - 2\gamma_k \lambda(f_{S_k}(x^k) - f_{S_k}^*) + 2\gamma_k \lambda(f_{S_k}(x^{k-1}) - f_{S_k}^*) \\ &\quad + \frac{4cL}{1+2cL} \gamma_k(f_{S_k}(x^k) - f_{S_k}^*) + \frac{2c^2L}{1+cL} \max\left\{\frac{2cL-1}{2cL+1}, 0\right\} f_{S_k}^* \\ &= \|z^k - x^*\|^2 - 2\rho(f_{S_k}(x^k) - f_{S_k}^*) - 2\left(\tilde{\gamma}_k + \gamma_k \lambda - \frac{2cL}{1+2cL} \gamma_k\right)(f_{S_k}(x^k) - f_{S_k}^*) \\ &\quad + 2\tilde{\gamma}_k(f_{S_k}(x^*) - f_{S_k}^*) + 2\gamma_k \lambda(f_{S_k}(x^{k-1}) - f_{S_k}^*) \\ &\quad + \frac{2c^2L}{1+cL} \max\left\{\frac{2cL-1}{2cL+1}, 0\right\} f_{S_k}^*. \end{aligned} \quad (23)$$

We need to find ρ such that

$$\tilde{\gamma}_k + \gamma_k \lambda - \frac{2cL}{1+2cL} \gamma_k \geq 0$$

Since $\tilde{\gamma}_k = \gamma_k - \rho$, then we have

$$\begin{aligned} \gamma_k - \rho + \gamma_k \lambda - \frac{2cL}{1+2cL} \gamma_k &\geq 0 \\ \Leftrightarrow \gamma_k \left(1 + \lambda - \frac{2cL}{1+2cL}\right) &\geq \rho. \end{aligned}$$

The inequality above is satisfied if it is satisfied for the lower bound on γ_k (which is $c/(1+cL)$), i.e.

$$\frac{c}{1+cL} \left(\frac{1}{1+2cL} + \lambda \right) \geq \rho.$$

We can take $\rho = \frac{c}{(1+cL)(1+2cL)}$ since $\lambda \geq 0$.

$$\begin{aligned} \tilde{\gamma}_k &= \gamma_k - \rho \\ &\leq c - \frac{c}{(1+cL)(1+2cL)} \\ &= c \frac{1+3cL+2c^2L^2-1}{(1+cL)(1+2cL)} \\ &\leq c^2L \frac{3+3cL}{(1+cL)(1+2cL)} \\ &= \frac{3c^2L}{1+2cL}. \end{aligned}$$

Using the above, we get from (23)

$$\begin{aligned} \|z^{k+1} - x^*\|^2 &\leq \|z^k - x^*\|^2 - 2\rho(f_{S_k}(x^k) - f_{S_k}(x^*)) + 2c\lambda(f_{S_k}(x^{k-1}) - f_{S_k}(x^*)) \\ &\quad + 2c\lambda(f_{S_k}(x^*) - f_{S_k}^*) + \frac{6c^2L}{1+2cL}(f_{S_k}(x^*) - f_{S_k}^*) \\ &\quad + \frac{2c^2L}{1+cL} \max \left\{ \frac{2cL-1}{2cL+1}, 0 \right\} f_{S_k}^*. \end{aligned}$$

Taking expectations we get

$$\begin{aligned} \mathbb{E} [\|z^{k+1} - x^*\|^2] &\leq \mathbb{E} [\|z^k - x^*\|^2] - 2\rho \mathbb{E} [f(x^k) - f(x^*)] + 2c\lambda \mathbb{E} [f(x^{k-1}) - f(x^*)] \\ &\quad + \left(2c\lambda + \frac{6c^2L}{1+2cL} \right) \sigma_{\text{int}}^2 + \frac{2c^2L}{1+cL} \max \left\{ \frac{2cL-1}{2cL+1}, 0 \right\} \sigma_{\text{pos}}^2. \end{aligned} \quad (24)$$

Rearranging terms we get

$$\begin{aligned} 2\rho \mathbb{E} [f(x^k) - f(x^*)] - 2c\lambda \mathbb{E} [f(x^{k-1}) - f(x^*)] &\leq \mathbb{E} [\|z^k - x^*\|^2] - \mathbb{E} [\|z^{k+1} - x^*\|^2] \\ &\quad + \left(2c\lambda + \frac{6c^2L}{1+2cL} \right) \sigma_{\text{int}}^2 \\ &\quad + \frac{2c^2L}{1+cL} \max \left\{ \frac{2cL-1}{2cL+1}, 0 \right\} \sigma_{\text{pos}}^2. \end{aligned} \quad (25)$$

Combining Equation (22) and Equation (25) for iterations $\{1, \dots, K-1\}$ we get

$$\begin{aligned} &2\rho \mathbb{E} [f(x^0) - f(x^*)] + 2\rho \sum_{k=1}^{K-1} \mathbb{E} [f(x^k) - f(x^*)] - 2c\lambda \sum_{k=1}^{K-1} \mathbb{E} [f(x^{k-1}) - f(x^*)] \\ &= 2\rho \sum_{k=0}^{K-1} \mathbb{E} [f(x^k) - f(x^*)] - 2c\lambda \sum_{k=0}^{K-2} \mathbb{E} [f(x^k) - f(x^*)] \\ &\leq (2\rho - 2c\lambda) \sum_{k=0}^{K-1} \mathbb{E} [f(x^k) - f(x^*)] \\ &\leq \|z^0 - x^*\|^2 + \frac{6c^2L}{1+2cL} \sigma_{\text{int}}^2 + \frac{2c^2L}{1+cL} \max \left\{ \frac{2cL-1}{2cL+1}, 0 \right\} \sigma_{\text{pos}}^2 \\ &\quad + \left(2c\lambda + \frac{6c^2L}{1+2cL} \right) (K-1) \sigma_{\text{int}}^2 + (K-1) \cdot \frac{2c^2L}{1+cL} \max \left\{ \frac{2cL-1}{2cL+1}, 0 \right\} \sigma_{\text{pos}}^2 \\ &\leq \|z^0 - x^*\|^2 + \left(2c\lambda + \frac{6c^2L}{1+2cL} \right) K \sigma_{\text{int}}^2 + K \cdot \frac{2c^2L}{1+cL} \max \left\{ \frac{2cL-1}{2cL+1}, 0 \right\} \sigma_{\text{pos}}^2. \end{aligned} \quad (26)$$

We need to ensure that $\rho - c\lambda > 0$ which is satisfied for λ such that

$$\begin{aligned}\frac{\rho}{2} &= \frac{c}{2(1+cL)(1+2cL)} > c\lambda \\ \Leftrightarrow 1 &> 2\lambda(1+cL)(1+2cL).\end{aligned}$$

Note that we also assume that $\lambda \leq cL$. Therefore, from (26) we get

$$\begin{aligned}\frac{1}{K} \sum_{k=0}^{K-1} \mathbb{E} [f(x^k) - f(x^*)] &\leq \frac{\|z^0 - x^*\|^2}{2(\rho - c\lambda)K} + \frac{1}{2(\rho - c\lambda)} \left(2c\lambda + \frac{6c^2L}{1+2cL} \right) \sigma_{\text{int}}^2 \\ &\quad + \frac{1}{2(\rho - c\lambda)} \frac{2c^2L}{1+cL} \max \left\{ \frac{2cL-1}{2cL+1}, 0 \right\} \sigma_{\text{pos}}^2 \\ &\leq \frac{\|z^0 - x^*\|^2}{2(\rho - c\lambda)K} + \frac{8c^2L}{2(\rho - c\lambda)} \sigma_{\text{int}}^2 \\ &\quad + \frac{1}{2(\rho - c\lambda)} \frac{2c^2L}{1+cL} \max \left\{ \frac{2cL-1}{2cL+1}, 0 \right\} \sigma_{\text{pos}}^2.\end{aligned}\tag{27}$$

Since $\rho - c\lambda \geq \frac{\rho}{2}$ and setting \bar{x}^k be uniformly at random chosen from $\{x^0, \dots, x^{K-1}\}$ we get

$$\mathbb{E} [f(\bar{x}^k) - f(x^*)] \leq \frac{\|z^0 - x^*\|^2}{\rho K} + \frac{8c^2L}{\rho} \sigma_{\text{int}}^2 + \frac{1}{\rho} \frac{2c^2L}{1+cL} \max \left\{ \frac{2cL-1}{2cL+1}, 0 \right\} \sigma_{\text{pos}}^2,\tag{28}$$

where we use the convexity of f and Jensen's inequality. Plugging the value of $\rho = \frac{c}{(1+cL)(1+2cL)}$ inside we get

$$\begin{aligned}\mathbb{E} [f(\bar{x}^k) - f(x^*)] &\leq \frac{\|z^0 - x^*\|^2}{cK} (1+cL)(1+2cL) + 8cL(1+cL)(1+2cL) \sigma_{\text{int}}^2 \\ &\quad + 2cL \max \{2cL-1, 0\} \sigma_{\text{pos}}^2.\end{aligned}\tag{29}$$

Choosing $c = \mathcal{O}(1/\sqrt{K})$ we get

$$\mathbb{E} [f(\bar{x}^k) - f(x^*)] \leq \mathcal{O} \left(\frac{\|z^0 - x^*\|^2}{\sqrt{K}} + \frac{\sigma_{\text{int}}^2}{\sqrt{K}} + \frac{\sigma_{\text{pos}}^2}{\sqrt{K}} \max \{2cL-1, 0\} \right).\tag{30}$$

Therefore, if $K \geq \mathcal{O}(\varepsilon^{-2})$ then $\mathbb{E} [f(\bar{x}^k) - f(x^*)] \leq \mathcal{O}(\varepsilon)$. It remains to notice that $z^0 = x^0$ to derive the statement of the theorem. \square

E Convergence of NGN-M with Decaying Step-size

Lemma 7. We have

$$\sum_{k=0}^{K-1} \frac{1}{k+1} \leq \log(K+2), \quad \sum_{k=0}^{K-1} \frac{1}{\sqrt{k+1}} \geq \frac{4}{5} \sqrt{K+1}.\tag{31}$$

Proof. We refer to Lemma A.8 from Garrigos and Gower [2023]. \square

To prove the convergence of NGN-M with decaying c_k we consider IMA formulation (see Section A in the paper):

$$\begin{aligned}x^{-1} &= z^0 = x^0, \quad z^{k+1} = z^k - \gamma_k \nabla f_{S_k}(x^k), \quad \gamma_k = \frac{c_k}{1 + \frac{c_k}{2f_{S_k}(x^k)} \|\nabla f_{S_k}(x^k)\|^2} \\ x^{k+1} &= \frac{\lambda}{1+\lambda} x^k + \frac{1}{1+\lambda} z^{k+1},\end{aligned}$$

where $c_k = \frac{c_0}{\sqrt{k+1}}$, $\lambda_k = Lc_k$, $\lambda_0 = 0$.

Theorem 4. Assume that each f_i is convex and L -smooth, and that Assumption 3.2 holds. Let the step-size hyperparameter is set $c_k = \frac{c_0}{\sqrt{k}}$, momentum parameter $\lambda_k \leq \min\{c_k L, 0.5(1 + c_k L)^{-1}(1 + 2c_k L)^{-1}\}$. Then the iterates of NGN-M satisfy

$$\begin{aligned} \mathbb{E} [f(\hat{x}^{K-1}) - f(x^*)] &\leq \frac{5(1 + c_0 L)(1 + 2c_0 L)\|x^0 - x^*\|^2}{4c_0 \sqrt{K}} + 10Lc_0(1 + c_0 L)(1 + 2c_0 L)\sigma_{\text{int}}^2 \frac{\log(K + 2)}{\sqrt{K}} \\ &\quad + 5c_0 L(1 + c_0 L) \frac{\log(K + 2)}{2\sqrt{K}} \max\{2c_0 L - 1, 0\} \sigma_{\text{pos}}^2, \end{aligned} \quad (32)$$

where $\hat{x}^{K-1} = \sum_{k=0}^{K-1} \frac{\rho_k}{\sum_{k=0}^{K-1} \rho_k} x^k$, $\rho_k = \frac{c_k}{(1 + c_k L)(1 + 2c_k L)}$.

Proof. At iteration $k = 0$ we have

$$z^1 = z^0 - \gamma_0 \nabla f_{S_0}(x^0) = x^0 - \gamma_0 \nabla f_{S_0}(x^0).$$

Therefore, we get

$$\begin{aligned} \|z^1 - x^*\|^2 &= \|z^0 - x^*\|^2 - 2\gamma_0 \langle \nabla f_{S_0}(x^0), z^0 - x^* \rangle + \gamma_0^2 \|\nabla f_{S_0}(x^0)\|^2 \\ &\stackrel{\text{Lem. B.6}}{\leq} \|z^0 - x^*\|^2 - 2\gamma_0 \langle \nabla f_{S_0}(x^0), x^0 - x^* \rangle + \frac{4c_0 L}{1 + 2c_0 L} \gamma_0 (f_{S_0}(x^0) - f_{S_0}^*) \\ &\quad + \frac{2c_0^2 L}{1 + c_0 L} \max\left\{\frac{2c_0 L - 1}{2c_0 L + 1}, 0\right\} f_{S_0}^*. \end{aligned} \quad (33)$$

Let $\gamma_0 = \rho_0 + \tilde{\gamma}_0$ where $\rho_0 = \frac{c_0}{(1 + c_0 L)(1 + 2c_0 L)}$. Then we have

$$\begin{aligned} \tilde{\gamma}_0 &= \gamma_0 - \rho_0 \\ &\stackrel{\text{Lem. B.5}}{\leq} c_0 - \frac{c_0}{(1 + c_0 L)(1 + 2c_0 L)} \\ &= c_0 \frac{1 + 3c_0 L + 2c_0^2 L^2 - 1}{(1 + c_0 L)(1 + 2c_0 L)} \\ &= c_0^2 L \frac{3 + 3c_0 L}{(1 + c_0 L)(1 + 2c_0 L)} \\ &= \frac{3c_0^2 L}{1 + 2c_0 L}. \end{aligned}$$

Using the above we continue from (33)

$$\begin{aligned} \|z^1 - x^*\|^2 &\stackrel{\text{conv.}}{\leq} \|z^0 - x^*\|^2 - 2\gamma_0 (f_{S_0}(x^0) - f_{S_0}(x^*)) + \frac{4c_0 L}{1 + 2c_0 L} \gamma_0 (f_{S_0}(x^0) - f_{S_0}^*) \\ &\quad + \frac{2c_0^2 L}{1 + c_0 L} \max\left\{\frac{2c_0 L - 1}{2c_0 L + 1}, 0\right\} f_{S_0}^* \\ &\leq \|z^0 - x^*\|^2 - 2\rho_0 (f_{S_0}(x^0) - f_{S_0}(x^*)) - 2\tilde{\gamma}_0 (f_{S_0}(x^0) - f_{S_0}^*) + 2\tilde{\gamma}_0 (f_{S_0}(x^*) - f_{S_0}^*) \\ &\quad + \frac{4c_0 L}{1 + 2c_0 L} \gamma_0 (f_{S_0}(x^0) - f_{S_0}^*) + \frac{2c_0^2 L}{1 + c_0 L} \max\left\{\frac{2c_0 L - 1}{2c_0 L + 1}, 0\right\} f_{S_0}^* \\ &= \|z^0 - x^*\|^2 - 2\rho_0 (f_{S_0}(x^0) - f_{S_0}(x^*)) - 2\left(\gamma_0 - \rho_0 - \frac{2c_0 L}{1 + 2c_0 L} \gamma_0\right) (f_{S_0}(x^0) - f_{S_0}^*) \\ &\quad + 2\tilde{\gamma}_0 (f_{S_0}(x^*) - f_{S_0}^*) + \frac{2c_0^2 L}{1 + c_0 L} \max\left\{\frac{2c_0 L - 1}{2c_0 L + 1}, 0\right\} f_{S_0}^*. \end{aligned} \quad (34)$$

Here we have

$$\begin{aligned} \gamma_0 - \rho_0 - \frac{2c_0 L}{1 + 2c_0 L} \gamma_0 &= \frac{1}{1 + 2c_0 L} \gamma_0 - \rho_0 \\ &= \frac{1}{1 + 2c_0 L} \gamma_0 - \frac{c_0}{(1 + c_0 L)(1 + 2c_0 L)} \\ &\stackrel{\text{Lem. B.5}}{\geq} \frac{1}{1 + 2c_0 L} \frac{c_0}{1 + c_0 L} - \frac{c_0}{(1 + c_0 L)(1 + 2c_0 L)} \\ &= 0, \end{aligned}$$

$\tilde{\gamma}_0 \leq \frac{3c_0^2 L}{1+2c_0 L}$, and $f_{S_0}(x^0) - f_{S_0}^* \geq 0$. Hence, we get

$$\begin{aligned} \|z^1 - x^*\|^2 &\leq \|z^0 - x^*\|^2 - 2\rho_0(f_{S_0}(x^0) - f_{S_0}^*) + \frac{6c_0^2 L}{1+2c_0 L}(f_{S_0}(x^*) - f_{S_0}^*) \\ &\quad + \frac{2c_0^2 L}{1+c_0 L} \max\left\{\frac{2c_0 L - 1}{2c_0 L + 1}, 0\right\} f_{S_0}^*. \end{aligned}$$

Rearranging terms and taking the expectation we get

$$\begin{aligned} 2\rho_0 \mathbb{E}[f(x^0) - f(x^*)] &\leq \mathbb{E}[\|z^1 - x^*\|^2] - \|z^0 - x^*\|^2 + \frac{6c_0^2 L}{1+2c_0 L} \sigma_{\text{int}}^2 \\ &\quad + \frac{2c_0^2 L}{1+c_0 L} \max\left\{\frac{2c_0 L - 1}{2c_0 L + 1}, 0\right\} \sigma_{\text{pos}}^2. \end{aligned} \quad (35)$$

Next, for $k > 0$ we can use the relation $z^k = x^k + \lambda_k(x^k - x^{k-1})$. We expand $\|z^{k+1} - x^*\|^2$

$$\begin{aligned} \|z^{k+1} - x^*\|^2 &= \|z^k - x^*\|^2 - 2\gamma_k \langle \nabla f_{S_k}(x^k), z^k - x^* \rangle + \gamma_k^2 \|\nabla f_{S_k}(x^k)\|^2 \\ &= \|z^k - x^*\|^2 - 2\gamma_k \langle \nabla f_{S_k}(x^k), x^k - x^* \rangle - 2\gamma_k \lambda_k \langle \nabla f_{S_k}(x^k), x^k - x^{k-1} \rangle \\ &\quad + \gamma_k^2 \|\nabla f_{S_k}(x^k)\|^2 \\ &\stackrel{\text{conv.}}{\leq} \|z^k - x^*\|^2 - 2\gamma_k(f_{S_k}(x^k) - f_{S_k}^*) - 2\gamma_k \lambda_k(f_{S_k}(x^k) - f_{S_k}(x^{k-1})) \\ &\quad + \gamma_k^2 \|\nabla f_{S_k}(x^k)\|^2 \\ &\stackrel{\text{Lem. B.6}}{\leq} \|z^k - x^*\|^2 - 2\gamma_k(f_{S_k}(x^k) - f_{S_k}^*) - 2\gamma_k \lambda_k(f_{S_k}(x^k) - f_{S_k}(x^{k-1})) \\ &\quad + \frac{4c_k L}{1+2c_k L} \gamma_k(f_{S_k}(x^k) - f_{S_k}^*) + \frac{2c_k^2 L}{1+c_k L} \max\left\{\frac{2c_k L - 1}{2c_k L + 1}, 0\right\} f_{S_k}^*. \end{aligned}$$

Let $\gamma_k = \rho_k + \tilde{\gamma}_k$, where $\rho, \tilde{\gamma}_k \geq 0$, and ρ is a constant step-size independent of S_k which will be defined later. Therefore, we have

$$\begin{aligned} \|z^{k+1} - x^*\|^2 &\leq \|z^k - x^*\|^2 - 2\rho_k(f_{S_k}(x^k) - f_{S_k}^*) - 2\tilde{\gamma}_k(f_{S_k}(x^k) - f_{S_k}^*) \\ &\quad - 2\gamma_k \lambda_k(f_{S_k}(x^k) - f_{S_k}^*) + 2\gamma_k \lambda_k(f_{S_k}(x^{k-1}) - f_{S_k}^*) \\ &\quad + \frac{4c_k L}{1+2c_k L} \gamma_k(f_{S_k}(x^k) - f_{S_k}^*) + \frac{2c_k^2 L}{1+c_k L} \max\left\{\frac{2c_k L - 1}{2c_k L + 1}, 0\right\} f_{S_k}^* \\ &= \|z^k - x^*\|^2 - 2\rho(f_{S_k}(x^k) - f_{S_k}^*) - 2\tilde{\gamma}_k(f_{S_k}(x^k) - f_{S_k}^*) + 2\tilde{\gamma}_k(f_{S_k}(x^*) - f_{S_k}^*) \\ &\quad - 2\gamma_k \lambda_k(f_{S_k}(x^k) - f_{S_k}^*) + 2\gamma_k \lambda_k(f_{S_k}(x^{k-1}) - f_{S_k}^*) \\ &\quad + \frac{4c_k L}{1+2c_k L} \gamma_k(f_{S_k}(x^k) - f_{S_k}^*) + \frac{2c_k^2 L}{1+c_k L} \max\left\{\frac{2c_k L - 1}{2c_k L + 1}, 0\right\} f_{S_k}^* \\ &= \|z^k - x^*\|^2 - 2\rho(f_{S_k}(x^k) - f_{S_k}^*) - 2\left(\tilde{\gamma}_k + \gamma_k \lambda - \frac{2c_k L}{1+2c_k L} \gamma_k\right)(f_{S_k}(x^k) - f_{S_k}^*) \\ &\quad + 2\tilde{\gamma}_k(f_{S_k}(x^*) - f_{S_k}^*) + 2\gamma_k \lambda_k(f_{S_k}(x^{k-1}) - f_{S_k}^*) \\ &\quad + \frac{2c_k^2 L}{1+c_k L} \max\left\{\frac{2c_k L - 1}{2c_k L + 1}, 0\right\} f_{S_k}^*. \end{aligned} \quad (36)$$

We need to find ρ_k such that

$$\tilde{\gamma}_k + \gamma_k \lambda - \frac{2c_k L}{1+2c_k L} \gamma_k \geq 0$$

Since $\tilde{\gamma}_k = \gamma_k - \rho_k$, then we have

$$\begin{aligned} \gamma_k - \rho_k + \gamma_k \lambda_k - \frac{2c_k L}{1+2c_k L} \gamma_k &\geq 0 \\ \Leftrightarrow \gamma_k \left(1 + \lambda_k - \frac{2c_k L}{1+2c_k L}\right) &\geq \rho_k. \end{aligned}$$

The inequality above is satisfied if it is satisfied for the lower bound on γ_k (which is $c/(1+cL)$), i.e.

$$\frac{c_k}{1+c_kL} \left(\frac{1}{1+2c_kL} + \lambda \right) \geq \rho.$$

We can take $\rho_k = \frac{c_k}{(1+c_kL)(1+2c_kL)}$ since $\lambda \geq 0$.

$$\begin{aligned} \tilde{\gamma}_k &= \gamma_k - \rho_k \\ &\leq c_k - \frac{c_k}{(1+c_kL)(1+2c_kL)} \\ &= c_k \frac{1+3c_kL+2c_k^2L^2-1}{(1+c_kL)(1+2c_kL)} \\ &\leq c_k^2L \frac{3+3c_kL}{(1+c_kL)(1+2c_kL)} \\ &= \frac{3c_k^2L}{1+2c_kL}. \end{aligned}$$

Using the above, we get from (36)

$$\begin{aligned} \|z^{k+1} - x^*\|^2 &\leq \|z^k - x^*\|^2 - 2\rho_k(f_{S_k}(x^k) - f_{S_k}(x^*)) + 2c_k\lambda_k(f_{S_k}(x^{k-1}) - f_{S_k}(x^*)) \\ &\quad + 2c_k\lambda_k(f_{S_k}(x^*) - f_{S_k}^*) + \frac{6c_k^2L}{1+2c_kL}(f_{S_k}(x^*) - f_{S_k}^*) \\ &\quad + \frac{2c_k^2L}{1+c_kL} \max \left\{ \frac{2c_kL-1}{2c_kL+1}, 0 \right\} f_{S_k}^*. \end{aligned}$$

Taking expectations, we get

$$\begin{aligned} \mathbb{E} [\|z^{k+1} - x^*\|^2] &\leq \mathbb{E} [\|z^k - x^*\|^2] - 2\rho_k \mathbb{E} [f(x^k) - f(x^*)] + 2c_k\lambda_k \mathbb{E} [f(x^{k-1}) - f(x^*)] \\ &\quad + \left(2c_k\lambda_k + \frac{6c_k^2L}{1+2c_kL} \right) \sigma_{\text{int}}^2 + \frac{2c_k^2L}{1+c_kL} \max \left\{ \frac{2c_kL-1}{2c_kL+1}, 0 \right\} \sigma_{\text{pos}}^2. \end{aligned} \quad (37)$$

Rearranging terms, we get

$$\begin{aligned} 2\rho_k \mathbb{E} [f(x^k) - f(x^*)] - 2c_k\lambda_k \mathbb{E} [f(x^{k-1}) - f(x^*)] &\leq \mathbb{E} [\|z^k - x^*\|^2] - \mathbb{E} [\|z^{k+1} - x^*\|^2] \\ &\quad + \left(2c_k\lambda_k + \frac{6c_k^2L}{1+2c_kL} \right) \sigma_{\text{int}}^2 \\ &\quad + \frac{2c_k^2L}{1+c_kL} \max \left\{ \frac{2c_kL-1}{2c_kL+1}, 0 \right\} \sigma_{\text{pos}}^2. \end{aligned} \quad (38)$$

Combining equation 35 and equation 38 for iterations $\{1, \dots, K-1\}$ we get

$$\begin{aligned} &2\rho_0 \mathbb{E} [f(x^0) - f(x^*)] + 2 \sum_{k=1}^{K-1} \rho_k \mathbb{E} [f(x^k) - f(x^*)] - 2 \sum_{k=1}^{K-1} c_k\lambda_k \mathbb{E} [f(x^{k-1}) - f(x^*)] \\ &= 2 \sum_{k=0}^{K-1} \rho_k \mathbb{E} [f(x^k) - f(x^*)] - 2 \sum_{k=0}^{K-2} c_k\lambda_k \mathbb{E} [f(x^k) - f(x^*)] \\ &\leq 2 \sum_{k=0}^{K-1} (\rho_k - c_k\lambda_k) \mathbb{E} [f(x^k) - f(x^*)] \\ &\leq \|z^0 - x^*\|^2 + \frac{6c_0^2L}{1+2c_0L} \sigma_{\text{int}}^2 + \frac{2c_0^2L}{1+c_0L} \max \left\{ \frac{2c_0L-1}{2c_0L+1}, 0 \right\} \sigma_{\text{pos}}^2 \\ &\quad + \sum_{k=1}^{K-1} \left(2c_k\lambda_k + \frac{6c_k^2L}{1+2c_kL} \right) \sigma_{\text{int}}^2 + \sum_{k=1}^{K-1} \frac{2c_k^2L}{1+c_kL} \max \left\{ \frac{2c_kL-1}{2c_kL+1}, 0 \right\} \sigma_{\text{pos}}^2. \end{aligned} \quad (39)$$

Note that choosing $\lambda_k = \min \{c_k L, 0.5(1 + c_k L)^{-1}(1 + 2c_k L)^{-1}\}$ ensures that $\frac{\rho_k}{2} \geq c_k \lambda_k$. Indeed, we have

$$\begin{aligned} \frac{\rho_k}{2} &= \frac{c_k}{2(1 + c_k L)(1 + 2c_k L)} > c_k \lambda_k \\ \Leftrightarrow 1 &> 2\lambda_k(1 + c_k L)(1 + 2c_k L). \end{aligned}$$

Therefore, from (39) and the facts that $\lambda_0 = 0$ and $\lambda_k \leq c_k L$ we get

$$\begin{aligned} \sum_{k=0}^{K-1} \rho_k \mathbb{E} [f(x^k) - f(x^*)] &\leq \|z^0 - x^*\|^2 + \sum_{k=0}^{K-1} \left(2c_k \lambda_k + \frac{6c_k^2 L}{1 + 2c_k L} \right) \sigma_{\text{int}}^2 \\ &\quad + \sum_{k=0}^{K-1} \frac{2c_k^2 L}{1 + c_k L} \max \left\{ \frac{2c_k L - 1}{2c_k L + 1}, 0 \right\} \sigma_{\text{pos}}^2 \\ &\leq \|z^0 - x^*\|^2 + 8L\sigma_{\text{int}}^2 \sum_{k=0}^{K-1} c_k^2 \\ &\quad + \sum_{k=0}^{K-1} 2c_k^2 L \max \left\{ \frac{2c_k L - 1}{2c_k L + 1}, 0 \right\} \sigma_{\text{pos}}^2. \end{aligned} \quad (40)$$

We have by Lemma 7

$$\begin{aligned} \sum_{k=0}^{K-1} \rho_k &= \sum_{k=0}^{K-1} \frac{c_k}{(1 + c_k L)(1 + 2c_k L)} \geq \sum_{k=0}^{K-1} \frac{c_k}{(1 + c_0 L)(1 + 2c_0 L)} \geq \frac{4c_0 \sqrt{K}}{5(1 + c_0 L)(1 + 2c_0 L)}, \\ \sum_{k=0}^{K-1} c_k^2 &\stackrel{\text{Lem 7}}{\leq} c_0^2 \log(K + 2), \\ \sum_{k=0}^{K-1} c_k^2 \max \left\{ \frac{2c_k L - 1}{2c_k L + 1}, 0 \right\} &\leq \sum_{k=0}^{K-1} c_k^2 \max \left\{ \frac{2c_0 L - 1}{2c_0 L + 1}, 0 \right\} \leq c_0^2 \log(K + 2) \max \left\{ \frac{2c_0 L - 1}{2c_0 L + 1}, 0 \right\}. \end{aligned} \quad (41)$$

Therefore, using equation 41, $z^0 = x^0$ in equation 40 and dividing both sides in equation 40 by $\sum_{k=0}^{K-1} \rho_k$ we derive

$$\begin{aligned} \sum_{k=0}^{K-1} \frac{\rho_k}{\sum_{k=0}^{K-1} \rho_k} \mathbb{E} [f(x^k) - f(x^*)] &\leq \frac{\|x^0 - x^*\|^2}{\sum_{k=0}^{K-1} \rho_k} + 8Lc_0^2 \sigma_{\text{int}}^2 \frac{\log(K + 2)}{\sum_{k=0}^{K-1} \rho_k} \\ &\quad + 2c_0^2 L \frac{\log(K + 2)}{\sum_{k=0}^{K-1} \rho_k} \max \left\{ \frac{2c_0 L - 1}{2c_0 L + 1}, 0 \right\} \sigma_{\text{pos}}^2. \end{aligned} \quad (42)$$

With an lower bound on $\sum_{k=0}^{K-1} \rho_k$ and Jensen's inequality we conclude that

$$\begin{aligned} \mathbb{E} [f(\hat{x}^{K-1}) - f(x^*)] &\leq \frac{5(1 + c_0 L)(1 + 2c_0 L) \|x^0 - x^*\|^2}{4c_0 \sqrt{K}} \\ &\quad + 10Lc_0(1 + c_0 L)(1 + 2c_0 L) \sigma_{\text{int}}^2 \frac{\log(K + 2)}{\sqrt{K}} \\ &\quad + 5c_0 L(1 + c_0 L) \frac{\log(K + 2)}{2\sqrt{K}} \max \{2c_0 L - 1, 0\} \sigma_{\text{pos}}^2, \end{aligned} \quad (43)$$

where $\hat{x}^{K-1} = \sum_{k=0}^{K-1} \frac{\rho_k}{\sum_{k=0}^{K-1} \rho_k} x^k$.

□

F Stability of NGN-M on a Simple Problem

We consider 1D convex functions of the form $f(x) = Lx^2(1 + p^2(x))$ that satisfy the following assumption.

Assumption 4. *There exists a constant C such that $C(1 + p^2(x)) \geq xp(x)p'(x)$.*

Note that $1 + p^2(x) \geq 1$ and $\deg(1 + p^2(x)) = \deg(xp(x)p'(x))$. Therefore, this assumption is mild.

Remark 3. For example, the function $f(x) = x^2(1 + x^2)$ (i.e., $p(x) = x$) is convex and satisfies Assumption 4 with $C = 1$.

Remark 4. Let $p(x) = \sum_{j=0}^m a_j x^j$. Then for large values of x in magnitude, $p(x) \sim a_m x^m$, $p'(x) \sim m a_m x^{m-1}$. Therefore, the constant C should be expected of order $C \approx m$, where $m = \deg(p(x))$.

The function $f(x)$ is non-negative for any $x \in \mathbb{R}$ and its minimum $f^* = 0$ is attained at $x = 0$ by design. Let us compute a step of NGN-M on this problem

$$\begin{aligned}
 x^{k+1} &= x^k - (1 - \beta) \frac{c}{1 + \frac{c}{2f(x^k)}(f'(x^k))^2} f'(x^k) + \beta(x^k - x^{k-1}) \\
 &= x^k - (1 - \beta) \frac{2Lc(1 + p^2(x^k) + x^k p(x^k)p'(x^k))}{1 + \frac{4L^2 c [x^k]^2}{2L[x^k]^2(1+p^2(x^k))}(1 + p^2(x^k) + x^k p(x^k)p'(x^k))^2} x^k + \beta(x^k - x^{k-1}) \\
 &= x^k - (1 - \beta) \underbrace{\frac{2Lc(1 + p^2(x^k) + x^k p(x^k)p'(x^k))}{1 + \frac{2Lc}{1+p^2(x^k)}(1 + p^2(x^k) + x^k p(x^k)p'(x^k))^2}}_{:=\hat{\gamma}_k} x^k + \beta(x^k - x^{k-1}). \tag{44}
 \end{aligned}$$

Note that the convexity of f implies that

$$\begin{aligned}
 f(0) &\geq f(x) + f'(x)(0 - x) \\
 0 &\geq Lx^2(1 + p^2(x)) - 2Lx^2(1 + p^2(x) + xp(x)p'(x)) \\
 0 &\geq -Lx^2(1 + p^2(x)) - 2Lx^3 p(x)p'(x) \\
 xp(x)p'(x) &\geq -\frac{1}{2}(1 + p^2(x)). \tag{45}
 \end{aligned}$$

In particular, equation 45 implies that $1 + p^2(x) + xp(x)p'(x) \geq \frac{1}{2}(1 + p^2(x)) > 0$. Therefore, we can obtain lower and upper bounds on $\hat{\gamma}_k$.

Lemma 8. Let Assumption 4 hold with a constant $C > 0$ and $f(x) = x^2(1 + p^2(x))$ be convex. Let $c \geq \frac{1}{2L}$. Then we have $\hat{\gamma}_k \in \left[\frac{1}{2(1+C)}, 2\right]$.

Proof. Indeed, the upper bound on $\hat{\gamma}_k$ follows from the following inequality

$$\begin{aligned}
 \hat{\gamma}_k &= \frac{2Lc(1 + p^2(x^k) + x^k p(x^k)p'(x^k))}{1 + \frac{2Lc}{1+p^2(x^k)}(1 + p^2(x^k) + x^k p(x^k)p'(x^k))^2} \\
 &\leq \frac{2Lc(1 + p^2(x^k) + x^k p(x^k)p'(x^k))}{\frac{2Lc}{1+p^2(x^k)}(1 + p^2(x^k) + x^k p(x^k)p'(x^k))^2} \\
 &= \frac{1 + p^2(x^k)}{1 + p^2(x^k) + x^k p(x^k)p'(x^k)} \leq 2, \tag{46}
 \end{aligned}$$

due to equation 45. The lower bound can be obtained as follows

$$\begin{aligned}
\hat{\gamma}_k &= \frac{2Lc(1 + p^2(x^k) + x^k p(x^k) p'(x^k))}{1 + \frac{2Lc}{1+p^2(x^k)}(1 + p^2(x^k) + x^k p(x^k) p'(x^k))^2} \\
&= \frac{2Lc(1 + p^2(x^k) + x^k p(x^k) p'(x^k))(1 + p^2(x^k))}{(1 + p^2(x^k)) + 2Lc(1 + p^2(x^k) + x^k p(x^k) p'(x^k))^2} \\
&\geq \frac{2Lc(1 + p^2(x^k) + x^k p(x^k) p'(x^k))(1 + p^2(x^k))}{2(1 + p^2(x^k) + x^k p(x^k) p'(x^k)) + 2Lc(1 + p^2(x^k) + x^k p(x^k) p'(x^k))^2} \\
&= \frac{Lc(1 + p^2(x^k))}{1 + Lc(1 + p^2(x^k) + x^k p(x^k) p'(x^k))} \\
&= \frac{Lc(1 + p^2(x^k))}{1 + Lc(1 + p^2(x^k) + C(1 + p^2(x^k)))} \\
&\geq \frac{Lc(1 + p^2(x^k))}{2Lc(1 + C)(1 + p^2(x^k))} = \frac{1}{2(1 + C)}
\end{aligned} \tag{47}$$

□

The update rule of NGN-M can be rewritten as

$$x^{k+1} = x^k - (1 - \beta)\hat{\gamma}_k x^k + \beta(x^k - x^{k-1}). \tag{48}$$

Let us consider the joint dynamics of $w^k := ([x^k]^\top, [x^{k-1}]^\top)^\top \in \mathbb{R}^{2d}$. We have that

$$\begin{aligned}
w^k &= \begin{pmatrix} x^k \\ x^{k-1} \end{pmatrix} = \begin{pmatrix} x^k - (1 - \beta)\hat{\gamma}_k x^k + \beta(x^k - x^{k-1}) \\ x^{k-1} \end{pmatrix} \\
&= \begin{pmatrix} \mathbf{I} - (1 - \beta)\hat{\gamma}_k \mathbf{I} + \beta \mathbf{I} & -\beta \mathbf{I} \\ \mathbf{I} & \mathbf{0} \end{pmatrix} \begin{pmatrix} x^k \\ x^{k-1} \end{pmatrix} = \mathbf{G} w^{k-1},
\end{aligned} \tag{49}$$

where

$$\mathbf{G} := \begin{pmatrix} \mathbf{I} - (1 - \beta)\hat{\gamma}_k \mathbf{I} + \beta \mathbf{I} & -\beta \mathbf{I} \\ \mathbf{I} & \mathbf{0} \end{pmatrix}. \tag{50}$$

Now we are ready to prove the convergence of NGN-M on this simple problem for any value $c \geq \frac{1}{2}$.

Theorem 5. *Let $f(x) = x^2(1 + p^2(x))$ be convex and Assumption 4 holds. Let $\beta \geq \frac{(2(1+C)-1)^2}{(2(1+C)+1)^2}$ and $c \geq \frac{1}{2L}$. Then the iterates of NGN-M on $f(x)$ converge to the minimum $f^* = 0$.*

Proof. We follow the standard proof of SGD with Polyak momentum [Polyak, 1964]. At this stage, we need to estimate the eigenvalues of \mathbf{G} . To do so, we will proceed with a permutation matrix Π ⁴ which transforms the matrix \mathbf{G} to the block-diagonal matrix as

$$\mathbf{G} = \begin{pmatrix} \mathbf{G}_1 & 0 & \dots & 0 \\ \dots & \dots & \dots & \dots \\ 0 & 0 & \dots & \mathbf{G}_d \end{pmatrix}, \tag{51}$$

where

$$\mathbf{G}_i := \begin{pmatrix} 1 + \beta - (1 - \beta)\hat{\gamma}_k & -\beta \\ 1 & 0 \end{pmatrix} \tag{52}$$

⁴The permutation matrix Π is defined as $\Pi_{ij} = \begin{cases} 1 & i \text{ odd}, j = i \\ 1 & i \text{ even}, j = 2n + i. \\ 0 & \text{else} \end{cases}$. Note that permutation matrices preserve eigenvalues.

Since the matrix \mathbf{G} is a block-diagonal matrix, we have $\|\mathbf{G}\| \leq \max_i \|\mathbf{G}_i\|$. Therefore, the problem is now simplified to bounding the spectral radii of the individual blocks \mathbf{G}_i , for $i = 1, 2, \dots, d$. The two eigenvalues u_1 and u_2 of \mathbf{G}_i are the roots of the quadratic

$$q(u) := u^2 - (1 + \beta - (1 - \beta)\hat{\gamma}_k)u + \beta = 0, \quad (53)$$

which take different values depending on the discriminant $\Delta := (1 + \beta - (1 - \beta)\hat{\gamma}_k)^2 - 4\beta$. Let us find the values of β when the discriminant is negative. We need to satisfy the inequality

$$\begin{aligned} (1 + \beta - (1 - \beta)\hat{\gamma}_k)^2 - 4\beta &\leq 0 \Leftrightarrow (1 + \beta)^2 + (1 - \beta)^2\hat{\gamma}_k^2 - 2(1 + \beta)(1 - \beta)\hat{\gamma}_k - 4\beta \leq 0 \\ &\Leftrightarrow (1 - \beta)^2 + (1 - \beta)^2\hat{\gamma}_k^2 - 2(1 + \beta)(1 - \beta)\hat{\gamma}_k \leq 0 \\ &\Leftrightarrow (1 - \beta)(1 + \hat{\gamma}_k^2) \leq 2(1 + \beta)\hat{\gamma}_k \\ &\Leftrightarrow \frac{1 + \hat{\gamma}_k^2}{2\hat{\gamma}_k} \leq \frac{1 + \beta}{1 - \beta}. \end{aligned} \quad (54)$$

Since the function $\frac{1+y^2}{2y}$ for $y \in \left[\frac{1}{2(1+C)}, 2\right]$ attains the maximum $\frac{4(1+C)^2+1}{4(1+C)}$ at $y = \frac{1}{2(1+C)}$, then we satisfy the last inequality, and consequently the discriminant is non-positive, if we choose

$$\frac{4(1+C)^2+1}{4(1+C)} \leq \frac{1+\beta}{1-\beta}. \quad (55)$$

The above inequality is satisfied for $\beta \in \left[\frac{(2(1+C)-1)^2}{(2(1+C)+1)^2}, 1\right)$. Therefore, we obtain that for such choice of β we have $\Delta_i \leq 0$ for all $i \in [d]$. Therefore, the zeros of the quadratic $q(u)$ are complex, and are equal in absolute value

$$|u_1| = |u_2| = \sqrt{\beta} < 1. \quad (56)$$

This gives us that $\|\mathbf{G}_i\| \leq \sqrt{\beta} < 1$. Therefore, the algorithm converges for any value of β in this range.

It remains to use Lemma 11 from [Foucart \[2012\]](#) which says that for a given matrix $\mathbf{A} \in \mathbb{R}^{d \times d}$, and $\epsilon > 0$, there exists a matrix norm $\|\cdot\|$ such that

$$\|\mathbf{A}\| \leq \rho(\mathbf{A}) + \epsilon, \quad (57)$$

where $\rho(\mathbf{A}) = \max\{|\lambda| : \lambda \text{ eigenvalue of } \mathbf{A}\}$ (spectral radius of \mathbf{A}).

Asymptotically ⁵ (as $k \rightarrow \infty$, one can show (see Theorem 12 in [Foucart \[2012\]](#)) that

$$\|w^k\|_2 = \mathcal{O}(\rho(\mathbf{G})^k), \quad (58)$$

where $\rho(\mathbf{G}) \leq \sqrt{\beta} < 1$ in our analysis. Therefore, NGN-M with hyperparameters $c \geq \frac{1}{2}$ and $\beta \geq \frac{1}{0}$ converges. \square

Remark 5. For example, NGN-M converges on $f(x) = x^2(1 + x^2)$ for any $c \geq \frac{1}{2}$ and $\beta \geq \frac{9}{25}$.

Theorem 5 shows that NGN-M remains stable even with an arbitrarily large step-size hyperparameter c . Thanks to the adaptive nature of NGN step-size, the actual update scale is automatically shrunk when necessary, preserving convergence. Importantly, this is possible with a choice of momentum parameter β close to 1, which extends the results of Section 4. We acknowledge that our current analysis is restricted to the special convex class of 1D functions $f(x) = x^2(1 + p^2(x))$ satisfying Assumption 4. Extending such stability guarantees to wider function classes with large momentum β remains a significant open challenge.

To support the theoretical result, we test the performance of NGN-M and GDM (Gradient Descent with Momentum) on the problem $f(x) = x^2(1 + x^2)$, which is convex and satisfies Assumption 4; see Figure F.1. We run both algorithms, varying the step-size hyperparameter in $\{10^{-4}, \dots, 10^4\}$. We

⁵A non-asymptotic version of the analysis can be derived using Theorem 5 by [Wang et al. \[2021\]](#)

run algorithms for 10^5 iterations. We stop training if the loss reaches a threshold 10^{-15} or exceeds 10^{10} for the first time. We observe that (i) for small step-size hyperparameters, both methods converge but do not reach the threshold 10^{-15} ; (ii) **NGN-M** reaches the threshold even for extremely large values of the step-size hyperparameter while **GDM** diverges. (iii) the fastest convergence of **GDM** is achieved with the step-size hyperparameter 10^{-2} after 691 iterations while the fastest convergence of **NGN-M** is achieved with $c = 10^1$ after 269 iterations. In other details, **NGN-M** achieves faster convergence and much more stable to the choice of the step-size hyperparameter. These results align well with our theoretical analysis.

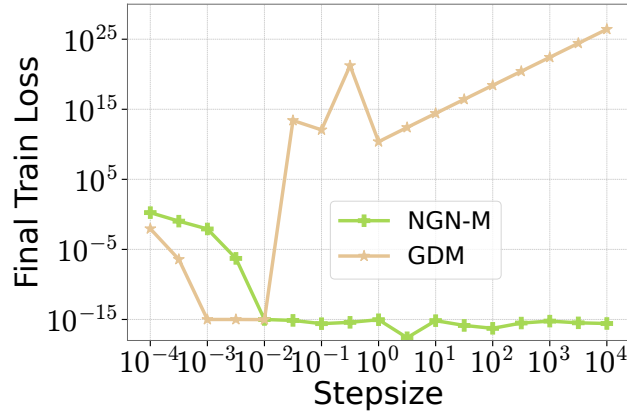


Figure F.1: Comparison of SGDM and NGN-M when minimizing a function $f(x) = x^2 + x^4$.

G How to Derive Diagonal NGN-based Step-size?

Here we provide derivations of how combine NGN and diagonal step-size following Section 3.3 for completeness.

We consider the following model

$$p^k = \arg \min_{p \in \mathbb{R}^d} \left[f_{\Sigma_k, c}(x^k + p) := (r(x^k) + \nabla r(x^k)^\top p)^2 + \frac{1}{2c} \|p\|_{\Sigma_k}^2 \right], \quad (59)$$

where $r(x) = \sqrt{f(x)}$. We compute the gradient of RHS of (59) w.r.t. p and equal it to zero:

$$\begin{aligned} \nabla_p f_{\Sigma_k, c}(x^k + p) &= 2 \left(r(x^k) + \nabla r(x^k)^\top p \right) \nabla r(x^k) + \frac{1}{c} \Sigma_k p \\ &= \left(2 \nabla r(x^k) \nabla r(x^k)^\top + \frac{1}{c} \Sigma_k \right) p + 2 r(x^k) \nabla r(x^k). \end{aligned}$$

Therefore, we have

$$p^k = - \left(2 \nabla r(x^k) \nabla r(x^k)^\top + \frac{1}{c} \Sigma_k \right)^{-1} 2 r(x^k) \nabla r(x^k).$$

Using Shermann-Morrison formula $(\mathbf{A} + uv^\top)^{-1} = \mathbf{A}^{-1} - \frac{\mathbf{A}^{-1}uv^\top \mathbf{A}^{-1}}{1 + u^\top \mathbf{A}^{-1}v}$ with $\mathbf{A} = 1/c \Sigma_k$ we derive

$$\begin{aligned} p^k &= - \left(c \Sigma_k^{-1} - \frac{2c^2 \Sigma_k^{-1} \nabla r(x^k) \nabla r(x^k)^\top \Sigma_k^{-1}}{1 + 2c \nabla r(x^k)^\top \Sigma_k^{-1} \nabla r(x^k)} \right) 2 r(x^k) \nabla r(x^k) \\ &= -2c r(x^k) \left(1 - \frac{2c \nabla r(x^k)^\top \Sigma_k^{-1} \nabla r(x^k)}{1 + 2c \nabla r(x^k)^\top \Sigma_k^{-1} \nabla r(x^k)} \right) \Sigma_k^{-1} \nabla r(x^k) \\ &= - \frac{2c r(x^k)}{1 + 2c \nabla r(x^k)^\top \Sigma_k^{-1} \nabla r(x^k)} \Sigma_k^{-1} \nabla r(x^k). \end{aligned}$$

Now we plug-in $r(x^k) = \sqrt{f(x^k)}$ and $\nabla r(x^k) = \frac{1}{2\sqrt{f(x^k)}} \nabla f(x^k)$ and obtain

$$\begin{aligned} p^k &= - \frac{2c \sqrt{f(x^k)}}{1 + 2c \frac{1}{4f(x^k)} \nabla f(x^k)^\top \Sigma_k^{-1} \nabla f(x^k)} \frac{1}{2\sqrt{f(x^k)}} \Sigma_k^{-1} \nabla f(x^k) \\ &= \frac{c}{1 + \frac{c}{2f(x^k)} \|\nabla f(x^k)\|_{\Sigma_k^{-1}}^2} \Sigma_k^{-1} \nabla f(x^k). \end{aligned}$$

G.1 Design Comparison of NGN-MDv1 and NGN-MDv2

The derivations in equation 3 are used to provide an intuition of how one can add a diagonal step-size into NGN by choosing the regularization matrix Σ_k . By choosing $\Sigma_k = \mathbf{D}_k$ we recover the update direction of NGN-MDv1. In this case, we have only one global NGN step-size in front of \mathbf{D}_k . The design of NGN-MDv2 follows a more straightforward intuition. In particular, it can be seen as a direct extension of NGN to diagonal case by replacing the squared gradient norm $\|\nabla f_{S_k}(x^k)\|^2$ by the squared partial derivative $(\nabla_j f_{S_k}(x^k))^2$ for each parameter $j \in [d]$.

The main difference in comparison with Adam is the order in which the preconditioning and momentum is applied. In both NGN-MDv1 and NGN-MDv2 we average the preconditioned updates $\Sigma_k^{-1} \nabla f_{S_k}(x^k)$, i.e. we first apply preconditioning and momentum later. In contrast, in Adam the stochastic gradients are averaged to construct new momentum term, and then the momentum is preconditioned. In other words, the momentum is applied first and then it is followed by preconditioning. We believe this change might be one of the reasons behind the step-size hyperparameter resilience as well.

In practice, we found out that the tuned performance of NGN-MDv1 is slightly better than that of NGN-MDv2. Moreover, NGN-MDv1 demonstrates higher resilience to the choice of the step-size hyperparameter than NGN-MDv2.

Table 2: Train time of Adam and NGN-MDv1 when training language models.

Model	Method	Time per Iteration (sec)	Time per Optimizer Update (sec)
70M	AdamW	1.63 ± 0.01	0.0048 ± 0.0002
	NGN-MDv1	1.65 ± 0.01	0.0130 ± 0.0002
160M	AdamW	3.33 ± 0.03	0.0088 ± 0.0003
	NGN-MDv1	3.37 ± 0.02	0.0239 ± 0.0003
410M	AdamW	8.41 ± 0.06	0.0838 ± 0.0009
	NGN-MDv1	8.68 ± 0.06	0.2154 ± 0.0007

G.2 Computation Cost of NGN-MD

Implementing any version of NGN-MD in practice might be slightly more computationally expensive. However, we highlight that computing a step of NGN-MD does not involve matrix-vector operations since the preconditioner is a diagonal matrix, and the matrix notation is used only for the convenience of presentation. The additional computation cost that we have in NGN-MDv1 is the computation of $\|\nabla f_{S_k}(x^k)\|_{\mathbf{D}_k^{-1}}^2$. This can naïvely be done by one additional pass over the gradient and summing the terms $\frac{1}{(\mathbf{D}_k)_{jj}}(\nabla_j f_{S_k}(x^k))^2$ for $j \in [d]$. This operation does not require additional matrix multiplication. However, it can be computed more efficiently while updating \mathbf{D}_k . The rest of the NGN-MDv1 implementation does not add any significantly costly operations in comparison with Adam.

We compare in Table 2 the time per iteration and optimizer update when training language models from Section 5 using AdamW and NGN-MDv1. We notice that our naive implementation of NGN-MDv1 is about 2.5 times slower than PyTorch’s AdamW. This is expected since our algorithm requires two passes over the gradient. Nevertheless, in this setting training time is dominated by forward and backward computations, keeping NGN-MDv1 competitive with AdamW. Moreover, as noted above, this overhead can be largely eliminated by computing the weighted gradient concurrently with the second-momentum v^k update. We do not aim to provide the most efficient implementation of NGN-MDv1 as the primary goal of our work is to highlight the stability advantages that NGN step-size brings in the training of neural networks.

G.2.1 Distributed Training

In a vanilla DDP implementation [Li et al., 2020], computing the weighted gradient norm $\|\nabla f_{S_k}(x^k)\|_{\mathbf{D}_k^{-1}}^2$ is straightforward since gradients are replicated across devices. We only require an additional all-reduce to synchronize $f_{S_k}(x^k)$ across devices, which is, however, a lightweight communication (just a single float) and, in principle, can even be overlapped with the backward pass.

However, with more sophisticated types of parallelism, like Tensor Parallel [Shoeybi et al., 2019] or ZeRO-2 [Rajbhandari et al., 2020], computing the weighted gradient norm introduces additional communication, as gradients are sharded across devices. This could still be implemented efficiently by accumulating squared gradient entries in each device and all-reducing only a single float, but it will, nevertheless, result in a computation and communication overhead for NGN-MDv1. We acknowledge that our methods might not be scalable to large distributed training, and adjustments are needed to make NGN-MDv1 work in this case. Nonetheless, we believe that our findings offer useful insights toward designing more stable optimization algorithms.

H How to add weight decay to NGN-MDv1?

Regularization techniques serve a fundamental purpose in minimizing generalization error. Orthogonal to their role for generalization, modern deep learning tasks often benefit from the use of weight decay [Xiao, 2024]. Despite its widespread application, the role of weight decay is poorly understood. Andriushchenko et al. [2023] suggested that it might provide implicit regularization by stabilizing the loss in over-parameterized neural networks and helping to balance the bias-variance tradeoff

that leads to lower training loss in under-parameterized networks. However, even in the case of SGD, there is still uncertainty regarding how the weight decay mechanism should be incorporated, as various implementations may exist [Zhang et al., 2018].

We propose two ways of adding weight decay to NGN-MDv1. The first variant follows the approach of Loshchilov and Hutter [2019], adding decoupled weight decay λ :

$$x^{k+1} = x^k - \lambda c x^k - (1 - \beta_1) \Sigma_k^{-1} \nabla f_{S_k}(x^k) + \beta_1(x^k - x^{k-1}). \quad (60)$$

In this update rule, the weight is added separately from the update direction $\Sigma_k^{-1} \nabla f_{S_k}(x^k)$. We call the resulting algorithm (60) Dec-NGN-MDv1, that stands for decoupled NGN-MDv1.

H.1 Combining NGN-MDv1 and Weight Decay Regularization

We now discuss how to combine NGN-MDv1 and weight decay, following the idea that weight decay should perform weight regularization.

We consider the following model

$$f_{\Sigma_k, \lambda}(x^k + p) := (r(x^k) + \nabla r(x^k)^\top p)^2 + \frac{1}{2c} \|p\|_{\Sigma_k}^2 + \frac{\lambda}{2} \|x^k + p\|_{\Sigma_k}^2.$$

By taking the gradient of $f_{\Sigma_k, \lambda}$ w.r.t. p we get

$$\begin{aligned} 0 &= 2(r(x^k) + \nabla r(x^k)^\top p) \nabla r(x^k) + \frac{1}{c} \Sigma_k p + \lambda \Sigma_k (x^k + p) \\ &= \left(2 \nabla r(x^k) \nabla r(x^k)^\top + \frac{1}{c} \Sigma_k + \lambda \Sigma_k \right) p + 2r(x^k) \nabla r(x^k) + \lambda \Sigma_k x^k. \end{aligned}$$

Therefore, we get

$$p^k = - \left(2 \nabla r(x^k) \nabla r(x^k)^\top + \frac{1}{c} \Sigma_k + \lambda \Sigma_k \right)^{-1} (2r(x^k) \nabla r(x^k) + \lambda \Sigma_k x^k).$$

Using Sherman-Morrison formula $(\mathbf{A} + uv^\top)^{-1} = \mathbf{A}^{-1} - \frac{\mathbf{A}^{-1}uv^\top \mathbf{A}^{-1}}{1 + u^\top \mathbf{A}^{-1}v}$ with $\mathbf{A} = (\lambda + 1/c) \Sigma_k$ and $u = v = \sqrt{2} \nabla r(x^k)$ we get that

$$\begin{aligned} &\left(2 \nabla r(x^k) \nabla r(x^k)^\top + \frac{1}{c} \Sigma_k + \lambda \Sigma_k \right)^{-1} \\ &= \frac{c}{1 + \lambda c} \Sigma_k^{-1} - \frac{\frac{2c^2}{(1+\lambda c)^2} \Sigma_k^{-1} \nabla r(x^k) \nabla r(x^k)^\top \Sigma_k^{-1}}{1 + \frac{2c}{1+\lambda c} \nabla r(x^k) \Sigma_k^{-1} \nabla r(x^k)}. \end{aligned}$$

Therefore, we have

$$\begin{aligned} p^k &= - \left(\frac{c}{1 + \lambda c} \Sigma_k^{-1} - \frac{\frac{2c^2}{(1+\lambda c)^2} \Sigma_k^{-1} \nabla r(x^k) \nabla r(x^k)^\top \Sigma_k^{-1}}{1 + \frac{2c}{1+\lambda c} \nabla r(x^k) \Sigma_k^{-1} \nabla r(x^k)} \right) (2r(x^k) \nabla r(x^k) + \lambda \Sigma_k x^k) \\ &= - \frac{2cr(x^k)}{1 + \lambda c} \left(1 - \frac{\frac{2c}{1+\lambda c} \nabla r(x^k)^\top \Sigma_k^{-1} \nabla r(x^k)}{1 + \frac{2c}{1+\lambda c} \nabla r(x^k) \Sigma_k^{-1} \nabla r(x^k)} \right) \Sigma_k \nabla r(x^k) \\ &\quad - \frac{\lambda c}{1 + \lambda c} x^k + \frac{\frac{2c^2 \lambda}{1+\lambda c} \Sigma_k^{-1} \nabla r(x^k) \nabla r(x^k)^\top x^k}{1 + \frac{2c}{1+\lambda c} \nabla r(x^k) \Sigma_k^{-1} \nabla r(x^k)} \\ &= - \frac{2cr(x^k)}{1 + \lambda c} \frac{1}{1 + \frac{2c}{1+\lambda c} \nabla r(x^k) \Sigma_k^{-1} \nabla r(x^k)} \Sigma_k^{-1} \nabla r(x^k) \\ &\quad - \frac{\lambda c}{1 + \lambda c} x^k + \frac{\frac{2c^2 \lambda}{1+\lambda c} \Sigma_k^{-1} \nabla r(x^k) \nabla r(x^k)^\top x^k}{1 + \frac{2c}{1+\lambda c} \nabla r(x^k) \Sigma_k^{-1} \nabla r(x^k)}. \end{aligned}$$

Algorithm 4 NGN-MDv1W

- 1: **Input:** $x^0 \in \mathbb{R}^d$, step-size parameter $c > 0$, momentum parameters $\beta_1, \beta_2 \in [0, 1]$, weight decay parameter $\lambda \geq 0$, stabilization parameter $\varepsilon > 0$
 - 2: **for** $k = 0, 1, \dots, K - 1$ **do**
 - 3: Sample a batch $S_k \subseteq [n]$ and compute f_{S_k} and $\nabla f_{S_k}(x^k)$
 - 4: Compute $v^k = \beta_2 v^{k-1} + (1 - \beta_2)(\nabla f_{S_k}(x^k) \odot \nabla f_{S_k}(x^k))$
 - 5: Compute $\mathbf{D}_k = \text{diag}(\varepsilon \mathbf{I} + \sqrt{v^k / (1 - \beta_2^k)})$
 - 6: Compute
$$\gamma_k = \frac{\frac{c}{(1+\lambda c)} \left[1 - \frac{c\lambda}{2f_{S_k}(x^k)} \nabla f_{S_k}(x^k)^\top x^k \right]_+}{1 + \frac{c}{2f_{S_k}(x^k)(1+\lambda c)} \|\nabla f_{S_k}(x^k)\|_{\mathbf{D}_k^{-1}}^2}$$
 - 7: Update $x^{k+1} = \frac{1}{1+\lambda c} x^k - (1 - \beta_1) \gamma_k \mathbf{D}_k^{-1} \nabla f_{S_k}(x^k) + \beta_1 (x^k - x^{k-1})$
 - 8: **end for**
- $[\cdot]_+$ denotes $\max\{0, \cdot\}$.
-

Using the connection $\nabla r(x^k) = \frac{1}{2\sqrt{f(x^k)}} \nabla f(x^k)$ and $r(x^k) = \sqrt{f(x^k)}$ we get

$$\begin{aligned} p^k &= -\frac{2c\sqrt{f(x^k)}}{1+\lambda c} \frac{1}{1 + \frac{2c}{4f(x^k)(1+\lambda c)} \nabla f(x^k)^\top \Sigma_k^{-1} \nabla f(x^k)} \Sigma_k^{-1} \frac{1}{2\sqrt{f(x^k)}} \nabla f(x^k) \\ &\quad - \frac{c\lambda}{1+\lambda c} x^k + \frac{\frac{2c^2\lambda}{4f(x^k)(1+\lambda c)} \Sigma_k^{-1} \nabla f(x^k) \nabla f(x^k)^\top x^k}{1 + \frac{2c}{4(1+\lambda c)f(x^k)} \nabla f(x^k)^\top \Sigma_k^{-1} \nabla f(x^k)} \\ &= -\frac{c/(1+\lambda c)}{1 + \frac{c}{2f(x^k)(1+\lambda c)} \|\nabla f(x^k)\|_{\Sigma_k^{-1}}^2} \Sigma_k \nabla f(x^k) - \frac{c\lambda}{1+\lambda c} x^k \\ &\quad + \frac{c\lambda}{1+\lambda c} \frac{\frac{c}{2f(x^k)} \nabla f(x^k)^\top x^k}{1 + \frac{c}{2f(x^k)(1+\lambda c)} \|\nabla f(x^k)\|_{\Sigma_k^{-1}}^2} \Sigma_k^{-1} \nabla f(x^k). \end{aligned}$$

To summarize, the update of NGN-Dv1W is the following

$$\begin{aligned} x^{k+1} &= x^k + p^k \\ &= \frac{1}{1+\lambda c} x^k + \frac{c\lambda}{1+\lambda c} \frac{\frac{c}{2f(x^k)} \nabla f(x^k)^\top x^k}{1 + \frac{c}{2f(x^k)(1+\lambda c)} \|\nabla f(x^k)\|_{\Sigma_k^{-1}}^2} \Sigma_k^{-1} \nabla f(x^k) \\ &\quad - \frac{c/(1+\lambda c)}{1 + \frac{c}{2f(x^k)(1+\lambda c)} \|\nabla f(x^k)\|_{\Sigma_k^{-1}}^2} \Sigma_k^{-1} \nabla f(x^k) \\ &= \frac{1}{1+\lambda c} x^k - \frac{\frac{c}{1+\lambda c} \left(1 - \frac{c\lambda}{2f(x^k)} \nabla f(x^k)^\top x^k \right)}{1 + \frac{c}{2f(x^k)(1+\lambda c)} \|\nabla f(x^k)\|_{\Sigma_k^{-1}}^2} \Sigma_k^{-1} \nabla f(x^k). \end{aligned} \tag{61}$$

To prevent the step-size next to $\Sigma_k^{-1} \nabla f(x^k)$ from being negative, the final update has the form

$$x^{k+1} = \frac{1}{1+\lambda c} x^k - \frac{\frac{c}{1+\lambda c} \left[1 - \frac{c\lambda}{2f(x^k)} \nabla f(x^k)^\top x^k \right]_+}{1 + \frac{c}{2f(x^k)(1+\lambda c)} \|\nabla f(x^k)\|_{\Sigma_k^{-1}}^2} \Sigma_k^{-1} \nabla f(x^k), \tag{62}$$

where $[\cdot]_+ := \max\{\cdot, 0\}$. Now we can add momentum on top and obtain the following update of NGN-MDv1W

$$x^{k+1} = \frac{1}{1+\lambda c} x^k - \frac{\frac{c}{1+\lambda c} \left[1 - \frac{c\lambda}{2f(x^k)} \nabla f(x^k)^\top x^k \right]_+}{1 + \frac{c}{2f(x^k)(1+\lambda c)} \|\nabla f(x^k)\|_{\Sigma_k^{-1}}^2} \Sigma_k^{-1} \nabla f(x^k) + \beta(x^k - x^{k-1}). \tag{63}$$

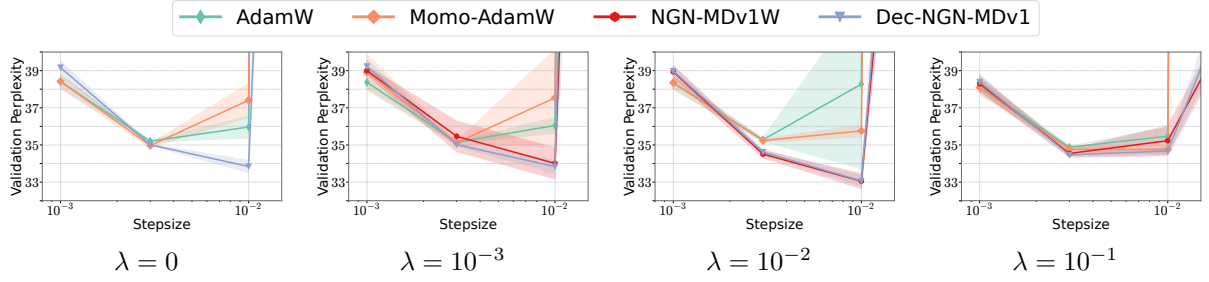


Figure H.1: Adding weight decay when pretraining a 70M Transformer++. When properly tuned, a value of weight decay > 0 enhances the performance of all algorithms. NGN-MDv1 retains his characteristic stability, and achieves smaller perplexity in all scenarios.

This combination of NGN-MDv1 and weight decay is summarized in Algorithm 4. We highlight that now the weight decay is incorporated inside the adaptive step-size as well as regularizing the coefficient next to x^k .

H.2 Empirical Validation of the Proposed Combinations

Having two possible ways of adding weight decay to NGN-MDv1, we test them on pretraining a 70M transformer on language modeling. The validation perplexity at the end of training is reported in Figure H.1. We note that when weight decay is turned off, both NGN-MDv1W and Dec-NGN-MDv1 reduce to NGN-MDv1.

First, we observe that when weight decay is properly tuned, all algorithms improve over the baseline case with no weight decay, which is consistent with the observation of Xiao [2024] and Andriushchenko et al. [2023] on AdamW. We also note that Dec-NGN-MDv1 and NGN-MDv1W require a smaller weight decay value compared to the other algorithms. Finally, the stability and performance of NGNMDv1 are preserved by both variations, allowing training with larger learning rates, and significantly improving over AdamW and Momo-Adam.

We do not observe a substantial difference between the two proposed modifications of NGN-MDv1 for this task. We remark however that these two versions serve substantially different purposes, and pretraining language models might not be the most representative task to evaluate the effect of adding regularization.

I Additional Experiments on Toy Problems

I.1 Additional Experiments on the Problem with Many Minima

Now, we provide a simple example of minimizing a function

$$f(x) = (\sin(1 + \cos(-\pi + x)) - 0.2x)^2 + (\sin(1 + \cos(\pi - x)) + 0.2x)^4 \quad (64)$$

that has many sharp sub-optimal local and flat global minima. We compare the performance of NGN-M and SGDM varying the step-size hyperparameter in $\{10^0, 10^1, 10^2, 10^3\}$ and the starting point in $[-20, 20]$ with a step $4/30$ ⁶. Based on the results in Figure I.1 (right), we conclude that (i) for small step-sizes, both methods likely get stuck at sub-optimal local minima and reach the global minima only if they are initialized close enough to it; (ii) for large step-sizes, we observe less runs of SGDM reaching the global minima; (iii) in contrast, for NGN-M with large step-sizes, we observe more runs reaching the global minima. This is possible due to the adaptive nature of the NGN step-size that forces NGN-M to converge to the flatness of the global minima.

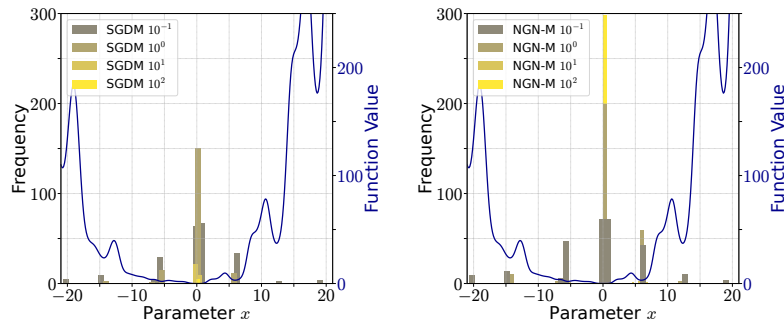


Figure I.1: Comparison of SGDM and NGN-M when minimizing function in equation 64.

I.2 Comparison on Rosenbrock Function

Now we present the results where we compare NGN-M and SGDM when minimizing the Rosenbrock function. We report the trajectories of optimizers and training dynamics in Figure I.2 and Figure I.3.

We observe that NGN-M converges for all values of c , indicating its high resilience to the choice of step-size hyperparameter. In contrast, SGDM already diverges for the step-size hyperparameter 10^{-2} . This can be explained by the adaptive nature of NGN step-size, which decreases the effective step-size of NGN-M for a more stable convergence. This is especially evident from the trajectories of algorithms. Indeed, NGN-M effectively moves in the complex valley of the Rosenbrock function, adapting to the local curvature.

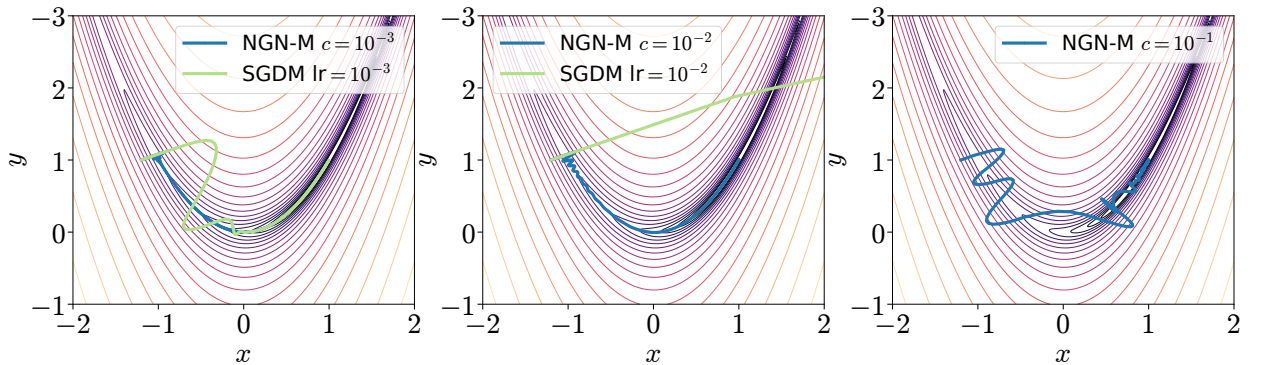


Figure I.2: Trajectories of NGN-M and SGDM when minimizing the Rosenbrock function and varying the step-size hyperparameter.

⁶This step is chosen small enough so that the initial point can be close to any local minima within $[-20, 20]$.

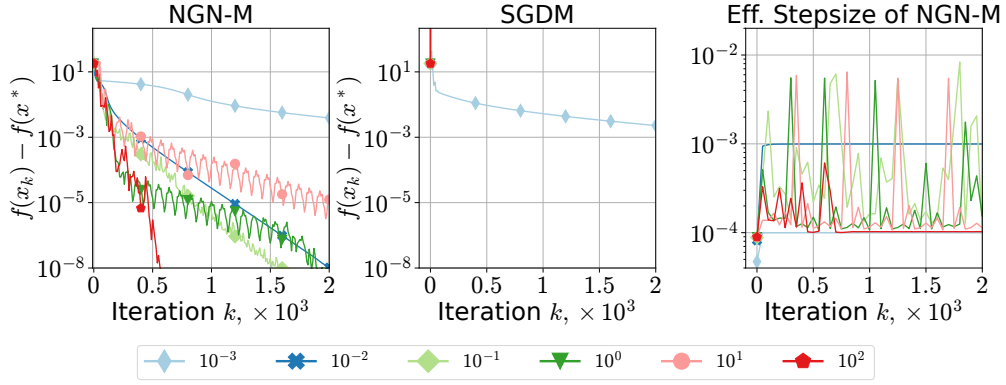


Figure I.3: Training dynamics of NGN-M and SGDM when minimizing the Rosenbrock function and varying the step-size hyperparameter.

I.3 Comparison on Quadratic Function with Theoretical Step-size

Next, we run NGN-M with theoretical choice of step-size hyperparameter $c = 1/\sqrt{K}$ and $c_k = 1/\sqrt{k}$ (see Theorem 1 and Theorem 4 for more details) against fixed choices $c \in \{10^{-3}, 10^{-4}\}$. The comparison is made on quadratic function $f(x) = \frac{1}{2}\|(\mathbf{A} + r\mathbf{I})x - y\|^2$, where $\mathbf{A} \in \mathbb{R}^{400 \times 400}$ and $y \in \mathbb{R}^{400}$ are sampled from standard normal distribution. The constant r controls the condition number of the problem.

We test the performance of NGN-M varying the condition number of the problem and the number of iterations; see Figure I.4. We observe that in all the cases, the choice $1/\sqrt{k}$ leads to faster convergence, supporting our theoretical claims. The choice $1/\sqrt{K}$ demonstrates competitive performance as well, but it is slightly pessimistic at the beginning of training. In contrast, the choice $c \in \{10^{-3}, 10^{-4}\}$, which is a default value in practice, is too small and does not lead to fast convergence.

These experiments demonstrate that when the problem satisfies all assumptions needed in the analysis, the choice of the step-size hyperparameter c given by the convergence theorems is a good starting point in practice and can serve as a baseline when tuning c .

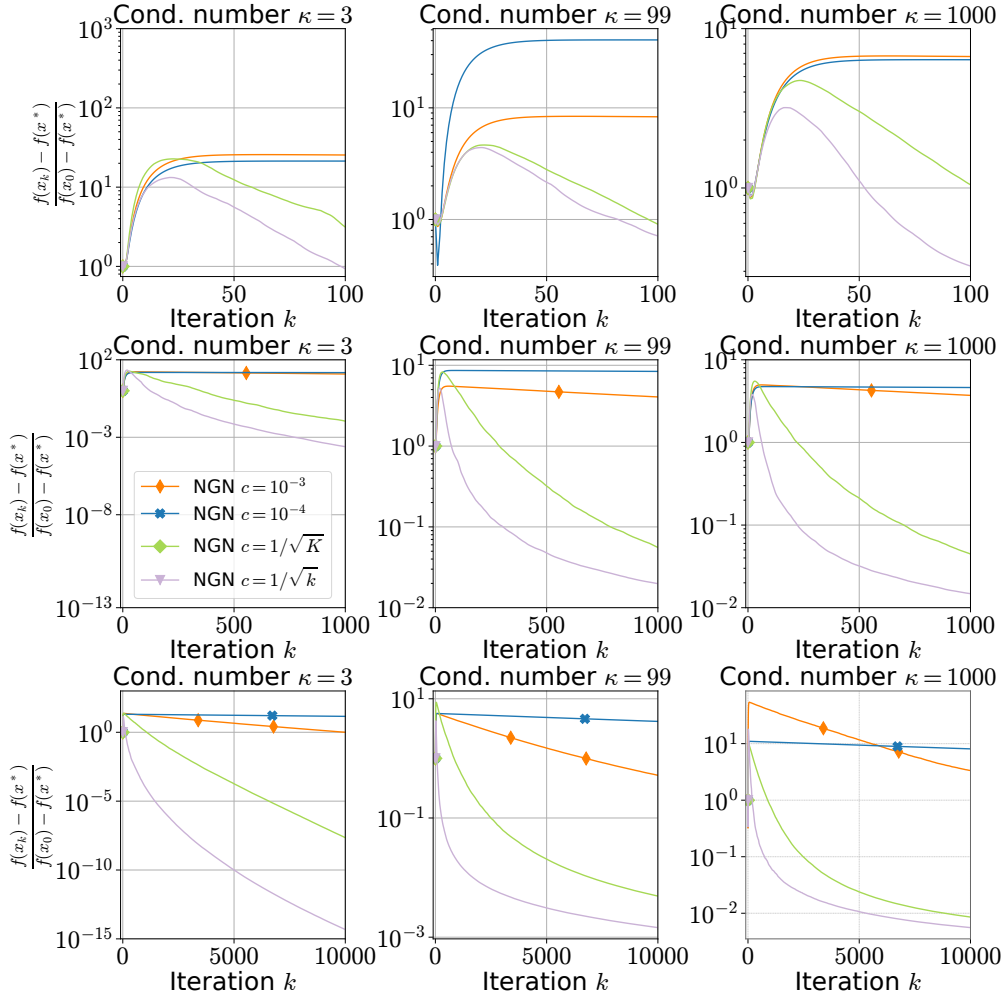


Figure I.4: Training dynamics of NGN-M with several choices of the step-size hyperparameter varying the condition number of the quadratic problem.

J Additional Experiments and Training Details

J.1 Training Details

The detailed experiment setup with hyperparameters and training details is presented in Table 3. We provide links to the exact model architectures used in our experiments (the links are clickable) as well as links to the tables and figures for each workload. We demonstrate the results averaged across 3 different random seeds for small and middle-range size experiments. We use standard values of momentum parameters $(\beta_1, \beta_2) = (0.9, 0.999)$ if the opposite is not specified. The step-size hyperparameter is tuned across powers of 10 (for some workloads we add additional values of the step-size hyperparameter shown in the step-size resilience plots). We use PyTorch [Paszke et al., 2017] implementation of Adam. The implementation of MomSPS, Momo, Momo-Adam are provided in the corresponding papers. Finally, when employing SGD-M, we set dampening equal to 0.9.

For vision transformers experiments, we follow the setup of Schaipp et al. [2024], and use Pytorch Image Models codebase [Wightman, 2019]. We train a vit_tiny_patch16_224 for 200 epochs on Imagenet1k, using a cosine learning rate schedule with a linear warmup of 5 epochs. Differently than Schaipp et al. [2024], we train in bfloat16, instead of float16, and do not employ weight decay regularization.

For pre-training Transformers on Causal Language Modeling, we build upon the nanoGPT [Karpathy, 2022] implementation, augmenting it with Rotational Positional Embedding [Su et al., 2023], RMSNorm [Zhang and Sennrich, 2019], and SwiGLU [Shazeer, 2020]. We call this enhanced version Transformer++. Models are trained with a batch size of 256, context length of 2048 tokens,

Table 3: Summary of experiment setup with all the details on hyperparameters used in each case.

Model	Dataset	Performance Results	Stability Results	Effective Stepsize Results	Epochs / Iterations	Batch Size	Comments
Resnet20	CIFAR10	Tab. 4, 5, 6	Fig. 2, J.1, J.2, J.5	Fig. J.9, J.10, J.6	50	128	
Resnet110	CIFAR100	Tab. 4, 5	Fig. 2, J.1, J.2, J.5		100	128	
VGG16	CIFAR10	Tab. 4, 5	Fig. J.1, J.2		50	128	
MLP	MNIST	Tab. 4, 5	Fig. J.1, J.3		10	128	2 hidden layers of size 100
ViT	CIFAR10	Tab. 4, 5	Fig. 2, J.1, J.2, J.5	Fig. 5, J.9, J.10, J.7	200	512	
LSTM	PTB	Tab. 5, 6	Fig. J.3		150	20	# layers 3
LSTM	Wiktext-2	Tab. 5, 6	Fig. J.8		150	20	# layers 3
Transformer	Rotten Tomatoes	Tab. 5, 6	Tab. J.8		2000	16	# heads 8 # layers 24
Transformer	Tiny Shakespeare	Tab. 5, 6	Fig. J.3, J.8		2000	16	# heads 8 # layers 24
Resnet18	ImageNet32	Tab. 4, 5,	Fig. J.4		45	128	constant learning rate schedule; no weight decay
Resnet18	ImageNet1k	Tab. 4, 5	Fig. 2, J.4		90	256	learning rate decay every 30 epochs by 0.1 no weight decay
ViT-Tiny	ImageNet1k	Tab. 5	Fig. 3		200	512	cosine learning rate schedule with linear warm-up for 5 epochs no weight decay, bfloat16
70M Transformer++	SlimPajama-627B	Tab. 5, 2	Fig. 4, H.1, J.14		2400	256	dim=512, # heads 8 # layers 6, context length 2048 (β_1, β_2) = (0.9, 0.95), bfloat16 clipping norm 1, linear warm-up for 10% of iterations
70M Transformer++	FineWeb	Tab. 7, 8, 9, 10			4800	128	dim=512, # heads 8 # layers 6, context length 2048 (β_1, β_2) = (0.9, 0.95), bfloat16 clipping norm 1, linear warm-up for 10% of iterations
160M Transformer++	SlimPajama-627B	Tab. 5, 2	Fig. 4, J.14	Fig. J.11, J.12, J.13	4800	256	dim=768, # heads 12 # layers 12, context length 2048 (β_1, β_2) = (0.9, 0.95), bfloat16 clipping norm 1, linear warm-up for 10% of iterations
410M Transformer++	SlimPajama-627B	Tab. 5, 2	Fig. 4, J.14		13500	256	dim=1024, # heads 16 # layers 24, context length 2048 (β_1, β_2) = (0.9, 0.95), bfloat16 clipping norm 1, linear warm-up for 10% of iterations
1B Transformer++	SlimPajama-627B	Tab. 5	Fig. 4, J.14		13500	256	dim=2048, # heads 8 # layers 16, context length 2048 (β_1, β_2) = (0.9, 0.95), bfloat16 clipping norm 1, linear warm-up for 10% of iterations

vocabulary size of 50280 and make use of GPT-Neox tokenizer [Black et al., 2022]. We adopt an enhanced training recipe, made popular by large language models such as LLaMa [Touvron et al., 2023]. These modifications include: training in **bfloat16**; employing a linear learning rate warm-up for 10% of the training steps, followed by cosine annealing to 10^{-5} ; omitting biases from linear layers; using $(\beta_1, \beta_2) = (0.9, 0.95)$ for all algorithms; clipping gradient norms above 1; no weight tying between embedding and last linear layer. All models are trained on SlimPajama-627B [Soboleva et al., 2023], a cleaned and deduplicated version of RedPajama. We report validation perplexity on a separate subset of Slim-Pajama consisting of 10M tokens. The total compute is estimated following Kaplan et al. [2020], where the estimated number of floating-point operations (FLOPs) is $6 \times \text{Number of Parameters} \times \text{Number of Tokens}$.

Experiments of small and middle size are performed on 1xRTX 4090. We perform ImageNet32 experiments on 2xA100-40GB, and ImageNet1k experiments on 4xA100-SXM4-40GB. For pretraining Transformers on Language Modeling, we employ 8xH100-HBM3-80GB GPUs. With multiple devices in use, we employ Distributed Data Parallel to parallelize the training process.

J.2 Comparison Algorithms that Support Momentum

In the main paper, we provided the test performance only. Now we additionally illustrate the performance of algorithms w.r.t. training loss convergence. Figure J.1 demonstrates that NGN-M is the most robust algorithm for the choice of the step-size hyperparameter from this perspective as well. In Figure J.1, we additionally demonstrate the performance of the algorithms on (VGG16 [Simonyan and Zisserman, 2014], CIFAR10) and (MLP, MNIST) workloads where NGN-M matches the performance of the state-of-the-art algorithms in this setting and archives higher resilience to the step-size hyperparameter choice. The best performance results are reported in Table 4 and

Table 4: The best validation score (with one standard deviation across 3 runs; accuracy for computer vision tasks; perplexity for NLP tasks) for the best learning rate choice for each method that supports momentum.

Model	Dataset	NGN	SGDM	NGN-M	MomSPS	Momo	ALR-SMAG
Resnet20	CIFAR10	88.30 \pm 0.20	85.42 \pm 0.70	88.76 \pm 0.05	87.20 \pm 0.38	88.86 \pm 0.14	88.88 \pm 0.19
Resnet110	CIFAR100	64.76 \pm 0.26	57.16 \pm 2.06	64.98 \pm 0.29	63.37 \pm 0.71	64.81 \pm 0.33	64.73 \pm 1.81
VGG16	CIFAR10	90.21 \pm 0.10	89.67 \pm 0.43	90.42 \pm 0.06	87.26 \pm 0.21	90.43 \pm 0.17	90.49 \pm 0.35
MLP	MNIST	98.04 \pm 0.07	97.63 \pm 0.10	97.97 \pm 0.08	97.73 \pm 0.09	97.97 \pm 0.04	97.64 \pm 0.06
ViT	CIFAR10	83.34 \pm 0.24	83.74 \pm 0.11	84.95 \pm 0.29	83.77 \pm 0.27	85.47 \pm 0.27	85.54 \pm 0.39
Resnet18	ImageNet32	48.63	48.56	48.29	N/A	48.68	N/A
Resnet18	ImageNet1k	67.00	66.73	67.12	N/A	67.09	N/A
Transformer	Tiny Shakespeare	9.27 \pm 0.19	8.73 \pm 0.13	7.67 \pm 0.12	N/A	8.80 \pm 0.19	N/A
Transformer	Rotten Tomatoes	9.01 \pm 0.22	8.75 \pm 0.04	7.12 \pm 0.03	N/A	8.65 \pm 0.03	N/A
LSTM	Wikitext-2	75.33 \pm 0.15	82.07 \pm 0.16	75.51 \pm 0.22	N/A	76.09 \pm 0.40	N/A

showcase that NGN-M always matches the performance of other optimizers or improves it.

J.3 Comparison of Algorithms that Support Momentum and Diagonal Step-size

Next, we illustrate the performance of the algorithms that support both momentum and diagonal step-size. According to the results in Figures J.2 and J.3, NGN-MDv1 achieves the best resilience to the step-size hyperparameter choice among all considered algorithms. Again, NGN-MDv1 is the most stable algorithm to the choice of step-size hyperparameter w.r.t. training loss convergence. Its best performance is competitive to that of other algorithms but the step-size hyperparameter range that gives such performance is wider.

Moreover, we support our claims about stability on additional workloads such as (VGG16, CIFAR10) (in Figure J.1), (MLP, MNIST), (LSTM [Hochreiter and Schmidhuber, 1997], PTB [Mikolov et al., 2010]), and (Transformer [Karpathy, 2022], Tiny Shakespeare [Karpathy, 2015]) workloads. We observe that NGN-MDv1 attains higher robustness to the choice of the step-size hyperparameter. Finally, the performance results on (LSTM, Wikitext-2 [Merity et al., 2016]) and (Transformer, Rotten Tomatoes [Pang and Lee, 2005]) are reported in Table 5. The results demonstrate competitive performance of NGN-MDv1 against other benchmarks across all considered workloads.

J.4 Additional ImageNet Experiments

Now we turn to the experiments involving training Resnet18 on ImageNet1k and ImageNet32. In Figure J.4 we provide the train loss curves and results on (Resnet18, ImageNet32) workload that demonstrate that NGN-M and NDN-MDv1 attain better resilience to the step-size hyperparameter choice than competitors not only from the train loss point of view as well. The best performance of algorithms is provided in Table 4 and 5. According to them, both NGN-M and NGN-M achieve competitive performance against considered benchmarks.

J.5 Additional Comparison against Lion, Adabelief, Adabound

This section compares algorithms from Section 5. Moreover, we include the comparison against Lion [Chen et al., 2024], Adabound [Luo et al., 2019], and Adabelief [Zhuang et al., 2020]. The results are presented in Table 5.

Table 5: The best validation score (with one standard deviation; accuracy for computer vision tasks; perplexity for NLP tasks) for the best learning rate choice for each method that supports diagonal step-sizes and momentum.

Model	Dataset	Adam	Momo-Adam	NGN-MDv1	NGN-MDv2	Lion	Adabelief	Adabound
Resnet20	CIFAR10	86.96 \pm 0.70	89.41 \pm 0.36	89.53 \pm 0.11	87.80 \pm 0.16	88.09 \pm 0.27	87.47 \pm 0.48	85.00 \pm 0.56
Resnet110	CIFAR100	64.12 \pm 0.94	67.10 \pm 0.53	66.10 \pm 0.45	64.33 \pm 0.40	61.85 \pm 0.77	65.32 \pm 0.43	61.28 \pm 0.39
VGG16	CIFAR10	90.26 \pm 0.23	90.95 \pm 0.28	90.64 \pm 0.18	90.07 \pm 0.37	N/A	N/A	N/A
MLP	MNIST	97.44 \pm 0.19	97.96 \pm 0.10	98.10 \pm 0.06	97.67 \pm 0.17	N/A	N/A	N/A
ViT	CIFAR10	85.96 \pm 0.23	85.74 \pm 0.12	85.65 \pm 0.10	86.56 \pm 0.11	86.89 \pm 0.19	85.05 \pm 0.47	80.32 \pm 0.47
Transformer	Rotten Tomatoes	6.80 \pm 0.07	6.81 \pm 0.05	6.90 \pm 0.05	6.83 \pm 0.05	N/A	N/A	N/A
Transformer	Tiny Shakespeare	6.80 \pm 0.06	6.80 \pm 0.05	6.89 \pm 0.06	6.82 \pm 0.05	N/A	N/A	N/A
LSTM	PTB	70.95 \pm 0.08	71.09 \pm 0.05	70.84 \pm 0.20	71.37 \pm 0.17	N/A	N/A	N/A
LSTM	Wikitext-2	81.49 \pm 1.49	82.23 \pm 0.64	75.24 \pm 0.21	81.99 \pm 0.78	N/A	N/A	N/A
Resnet18	ImageNet32	48.11	48.09	48.06	47.55	N/A	N/A	N/A
Resnet18	ImageNet1k	67.17	67.06	67.15	67.32	N/A	N/A	N/A
ViT-Tiny	ImageNet1k	71.05 \pm 0.16	71.22 \pm 0.36	71.345 \pm 0.22	N/A	N/A	N/A	N/A
Transformer++ 70M	SlimPajama-627B	34.38 \pm 0.12	34.96 \pm 0.11	33.84 \pm 0.33	N/A	N/A	N/A	N/A
Transformer++ 160M	SlimPajama-627B	24.03 \pm 0.02	24.29 \pm 0.10	23.32 \pm 0.06	N/A	N/A	N/A	N/A
Transformer++ 410M	SlimPajama-627B	16.65 \pm 0.03	17.07 \pm 0.05	16.48 \pm 0.03	N/A	N/A	N/A	N/A
Transformer++ 1B	SlimPajama-627B	13.09	N/A	13.11	N/A	N/A	N/A	N/A

We observe that NGN-MDv1 and NGN-MDv2 both achieve competitive performance across various Deep Learning workloads. In Figure J.5, we observe that Lion, Adabound and Adabelief algorithms do not match always the performance of NGN-MDv1 and Adam: Adabelief has worse performance on (Resnet20, CIFAR10) workload; Adabound has worse performance on (Resnet20, CIFAR10), (Resnet110, CIFAR100), and (ViT, CIFAR10) workloads; Lion has worse performance on (Resnet110, CIFAR100) workload. Moreover, their resilience to the step-size hyperparameter choice is lower than that of NGN-MDv1. To summarize, NGN-M and NGN-MDv1 are the most robust algorithms to the choice of step-size hyperparameter.

J.6 Comparison of Adaptive Step-sizes of Adam, Momo-Adam, and NGN-MDv1

Next, we conduct experiments to compare the adaptive step-size of Adam, Momo-Adam, and NGN-MDv1. Note that ResNet20 model consists of 3 base blocks, and each block has 3 convolution layers. In Figure J.6 we plot the average adaptive step-size of the layers $j \in \{\text{layer1.0.conv1, layer2.0.conv1, layer3.0.conv1}\}$ of ResNet20 that corresponds to the first convolution layer within each base block. Similarly, in Figure J.7 we plot the average adaptive step-size of the layers $j \in \{\text{layer0.0.fn.to_qkv, layer3.0.fn.to_qkv, layer5.0.fn.to_qkv}\}$ that corresponds to the attention layers of the first, fourth, and sixth base blocks.

Since the adaptivity of Adam is only in the second-order momentum applied as a normalization, in our experiment we compare the following quantities

$$\frac{\gamma}{(\mathbf{D}_k)_{(j)}} \text{ for Adam, } \frac{\tau_k}{(\mathbf{D}_k)_{(j)}} \text{ for Momo-Adam, } \frac{\gamma_k}{(\mathbf{D}_k)_{(j)}} \text{ for NGN-MDv1,} \quad (65)$$

where γ is the step-size hyperparameter of Adam.

Let us first describe the results for ResNet20 in Figure J.6. We observe that NGN-MDv1 tends to set smaller effective step-size compared to two other algorithms. This is especially visible for the large step-size hyperparameter values where the adaptive step-size of NGN-MDv1 is by several orders in magnitude smaller than that of Adam and Momo-Adam. In contrast, the coordinate-wise

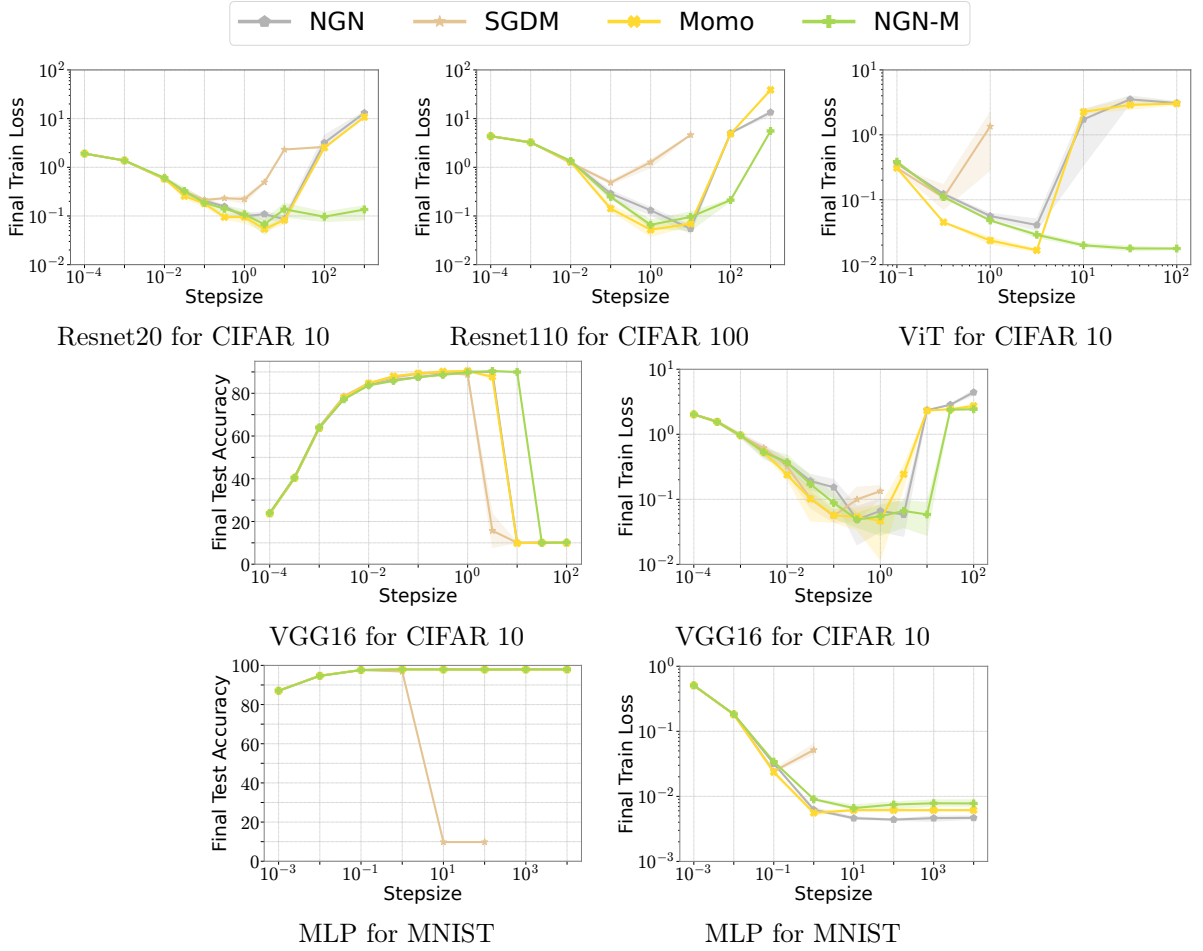


Figure J.1: Stability performance of algorithms supporting momentum varying step-size hyperparameter (c for NGN and NGN-M, α_0 for Momo, and step-size for SGDM). We observe that NGN-M achieves the training loss close to the best possible for a wider range of the step-size hyperparameter.

adaptive step-size of Momo-Adam is mostly follow that of Adam. Considering that the stability performance of NGN-MDv1 is much higher for this task, this happens mainly due to the fact that the adaptation mechanism of NGN-MDv1 step-size is more conservative than that of Momo-Adam.

Now we switch to the results on ViT model in Figure J.7. Here both Momo-Adam and NGN-MDv1 tend to utilize smaller effective coordinate-wise step-size, by several orders in magnitude smaller than that of Adam. However, the adaptation mechanism of NGN-MDv1 is still more conservative than that of Momo-Adam, especially for large step-size hyperparameters. We also highlight that in this experiment the best performance of NGN-MDv1 is achieved with $c = 10^{-3}$. When we vary the step-size hyperparameter c , the effective coordinate-wise step-size does not change dramatically, especially for layers.0.0.fin.to_qkv layer.

J.7 Extended Comparison of Momentum-based Algorithms on NLP Tasks

We switch to comparison of NGN-M, Momo, NGN, and SGDM on NLP tasks. In particular, we consider the training of Transformer (based on NanoGPT) on the Tiny Shakespeare and Rotten Tomatoes datasets and LSTM on the Wikitext-2 dataset from Appendix J.3. We report the results in Figure J.8 while the best performance is shown in Table 4. First, note that all algorithms do not match the best performance of those that incorporate diagonal step-size and momentum (see Table 5). Such results are expected since the training of NLP models has significantly different coordinate-wise conditioning. Nonetheless, NGN-M algorithm achieves better resilience to the step-size hyperparameter choice, especially in the training of Transformer models. Therefore, NGN-M across various model architectures and task domains.

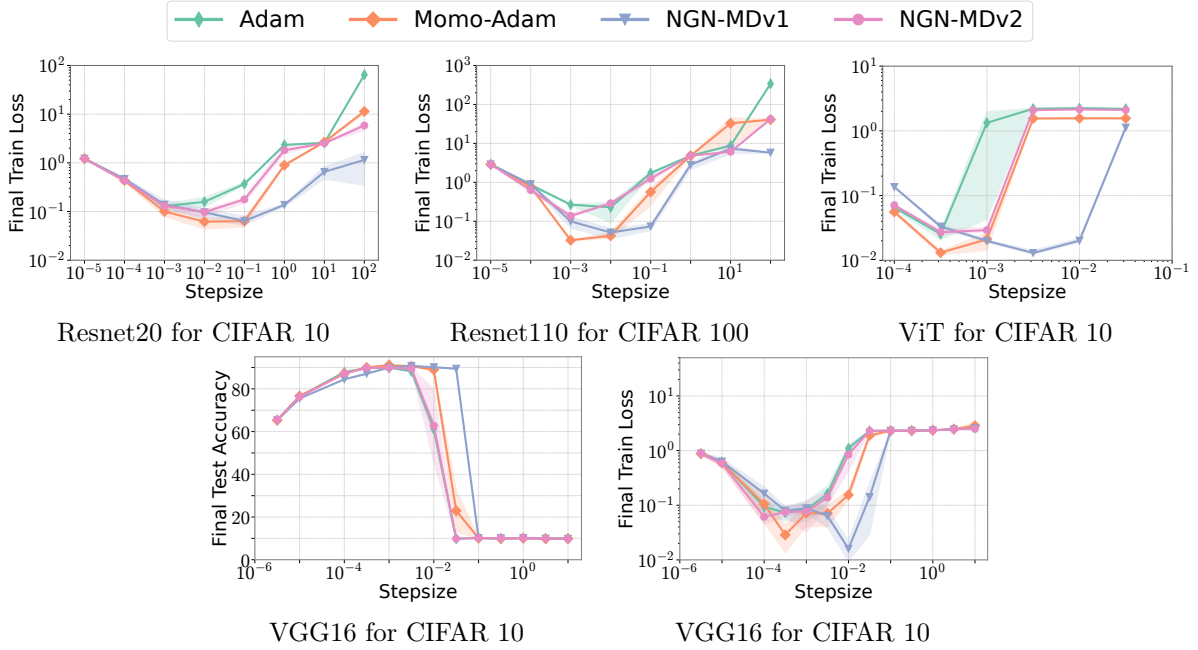


Figure J.2: Stability performance of algorithms supporting momentum and diagonal step-size varying step-size hyperparameter (c for NGN-MDv1 and NGN-MDv2, α_0 for Momo-Adam, and step-size for Adam). We observe that NGN-MDv1 achieves the training loss close to the best possible for a wider range of the step-size hyperparameter.

Table 6: The best validation score (with one standard deviation; accuracy for image classification; perplexity for language modeling) for the best learning rate choice for each method that supports diagonal step-sizes.

Model	Dataset	Adagrad	RMSprop	NGN-D
Resnet20	CIFAR10	85.90 \pm 0.30	86.71 \pm 0.64	86.98 \pm 0.15
Transformer	Rotten Tomatoes	7.77 \pm 0.02	6.87 \pm 0.05	6.92 \pm 0.03
Transformer	Tiny Sheaksper	7.77 \pm 0.05	7.00 \pm 0.13	6.90 \pm 0.05
LSTM	PTB	99.24 \pm 2.13	69.00 \pm 0.17	71.54 \pm 0.11
LSTM	Wikitext-2	113.19 \pm 4.36	79.48 \pm 0.45	75.44 \pm 0.12

J.8 Comparison of Algorithms with Diagonal Step-size

Now we compare algorithms with diagonal step-size such as NGN-D, Adagrad [Duchi et al. \[2011\]](#), and RMSprop [Kingma and Ba \[2015\]](#). Since NGN-D requires to find constants $\{c_j\}_{j=1}^d$ where d is the size of the model. Finding sufficiently good constants c_j might be a challenging task since d is a large number. Therefore, we use RMSprop preconditioner \mathbf{D}_k to set them as $c_j = c/(\mathbf{D}_k)_{(j)}$. We leave the exploration of how to set constants c_j properly for future research.

For each method, we tune its learning rate hyperparameter over the powers of 10: $\{10^{-4}, \dots, 10^2\}$ and present the best performance averaged across 3 random seeds in Table 6. We observe that NGN-D performs similarly to RMSprop. NGN-D has slightly worse performance on (LSTM, PTB) dataset but significantly better on (LSTM, Wikitext-2) workload. Besides, Adagrad always has the worst performance. Moreover, these algorithms do not have high resilience to the choice of hyperparameter. Therefore, we omit their comparison from this perspective.

J.9 Effective Step-size of NGN-M, Momo, NGN-MDv1, and Momo-Adam

Next, we compare the effective step-size applied throughout the training with NGN-M, Momo, NGN-MDv1, and Momo-Adam in Figures J.9 and J.10. First, both NGN-M and Momo perform a warm-up

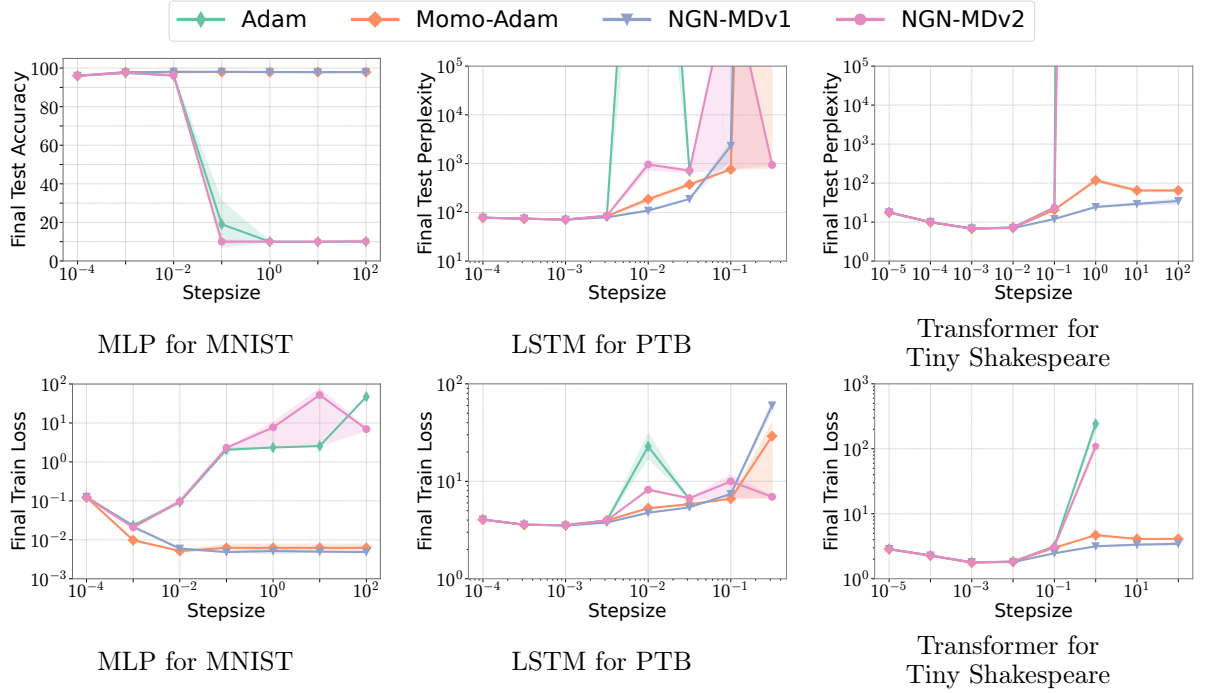


Figure J.3: Stability performance of algorithms supporting momentum and diagonal step-size varying step-size hyperparameter (c for NGN-MDv1 and NGN-MDv2, α_0 for Momo-Adam, and step-size for Adam). We observe that NGN-MDv1 achieves the training loss close to the best possible for a wider range of the step-size hyperparameter.

in the beginning: the effective step-size increases at the beginning of the training. Then we observe the main difference between the two algorithms above: effective step-size of Momo for sufficiently large step-size hyperparameter is not adaptive within some part of the training, it always hits the upper bound. Consequently, during that part of the training Momo reduces to SGDM. In contrast, the effective step-size of NGN-M is always adaptive: it gradually decreases after a short warm-up. This trend is similar to the state-of-the-art learning rate schedulers used in practice. Similar observations can be made in comparison of NGN-MDv1 and Momo-Adam.

J.10 Effective Updates in Training Language Models

In this section, we demonstrate the magnitude of updates when training 160M language model with Adam and NGN-MDv1 and varying the step-size hyperparameter across different layers of the model: see the results in Figures J.11 to J.13. We demonstrate that NGN-MDv1 is a more conservative algorithm: the effective update is smaller than that of Adam due to the adaptive nature of the step-size. This is especially evident when training 160M language model with a step-size hyperparameter 0.03: The updates of Adam become considerably larger than the update of NGN-MDv1. This property is a key factor behind the difference in training dynamics: NGN-MDv1 can stabilize at a significantly lower training loss.

J.11 Training Dynamics in Training Language Models

Now we report the training dynamics in the training language across all tested sizes.

J.12 Ablation Study of Momentum Parameters

In this section, we study the sensitivity of NGN-MDv1 and Adam to the choice of the learning rate and momentum hyperparameters, when training 70M language model on FineWeb dataset [Penedo et al., 2024]. To do that, we fix $\beta_1 = 0.9$ (or $\beta_2 = 0.95$) and make a sweep over the learning

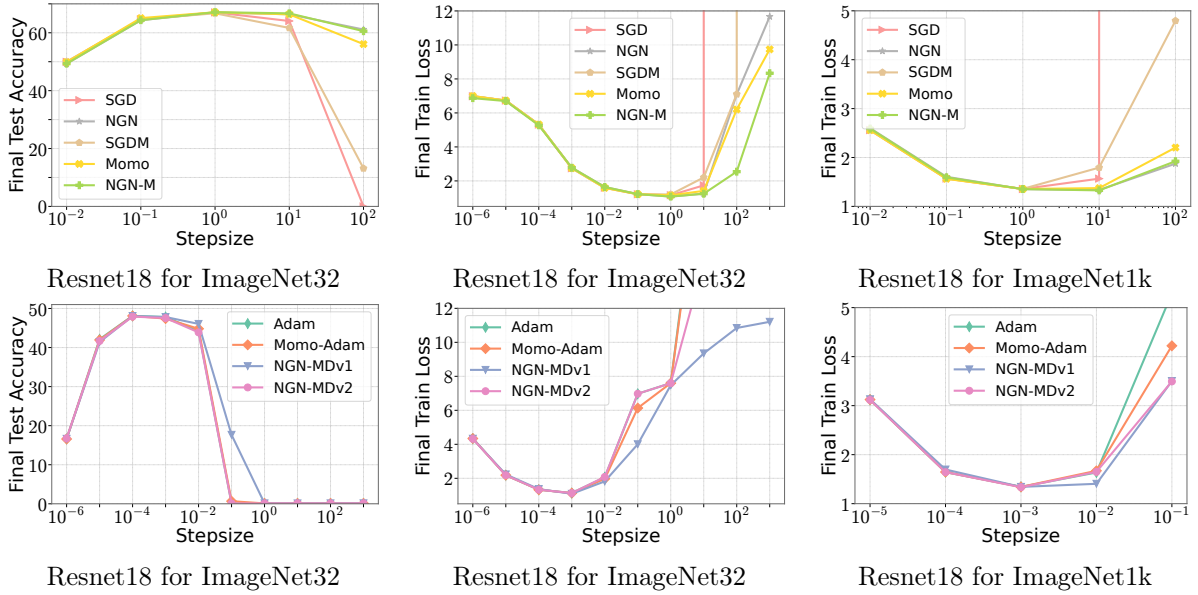


Figure J.4: Stability performance of algorithms supporting momentum (**first row**), and momentum with diagonal step-size (**second row**) varying step-size hyperparameter (c for NGN, NGN-M, NGN-MDv1, and NGN-MDv2, α_0 for Momo and Momo-Adam, and step-size for SGD, SGDM, and Adam).

Table 7: Test perplexity of NGN-MDv1 when varying the learning rate and β_1 hyperparameters when training 70M language model on the FineWeb dataset.

lr	$\beta_1 = 0.6$	$\beta_1 = 0.8$	$\beta_1 = 0.9$	$\beta_1 = 0.99$
3e-4	49.9 \pm 0.2	47.4 \pm 0.2	47.0 \pm 0.2	49.7 \pm 0.3
1e-3	41.5 \pm 0.2	39.9 \pm 0.2	38.6 \pm 0.1	40.2 \pm 0.3
3e-3	40 \pm 1	36.9 \pm 0.3	35.9 \pm 0.1	37.2 \pm 0.1
1e-2	54 \pm 16	37 \pm 2	34.7 \pm 0.3	35.9 \pm 0.1
3e-2	278 \pm 6	129 \pm 2	34.6 \pm 0.1	35.6 \pm 0.1

Table 8: Test perplexity of Adam when varying the learning rate and β_1 hyperparameters when training 70M language model on the FineWeb dataset.

lr	$\beta_1 = 0.6$	$\beta_1 = 0.8$	$\beta_1 = 0.9$	$\beta_1 = 0.99$
3e-4	49.4 \pm 0.2	46.5 \pm 0.1	46.2 \pm 0.3	57 \pm 1
1e-3	41.4 \pm 0.2	39.6 \pm 0.1	38.5 \pm 0.1	45.0 \pm 0.2
3e-3	40.7 \pm 0.1	37.0 \pm 0.1	36.0 \pm 0.1	220 \pm 70
1e-2	160 \pm 60	41 \pm 2	36 \pm 2	210 \pm 110
3e-2	420 \pm 20	340 \pm 50	320 \pm 60	330 \pm 130

rate hyperparameter and β_2 (or learning rate hyperparameter and β_1). We report the final test perplexity averaged over 3 runs for each set of hyperparameters.

We summarize our findings from Table 7, Table 8, Table 9, and Table 10 as follows:

- Low lr (3e-3): NGN-MDv1 and Adam show similar sensitivity to changes in both β_1 and β_2 .
- Moderate lr (1e-2): NGN-MDv1 is noticeably more robust than Adam to extremes of β_1 , while both optimizers perform similarly across β_2 (though Adam’s performance degrades slightly at $\beta_2 = 0.999$).
- High lr (3e-2): Both methods suffer when β_1 is small (or β_2 is large), but NGN-MDv1 recovers lower perplexity at larger β_1 values (smaller β_2 values), whereas Adam fails to reach comparable performance.

To conclude, NGN-MDv1 demonstrates greater robustness to changes in momentum parameters at high lr, and consistently attains lower perplexity than Adam, even when both methods’ performance deteriorates (we refer to the cases when both algorithms cannot achieve perplexity around 50).

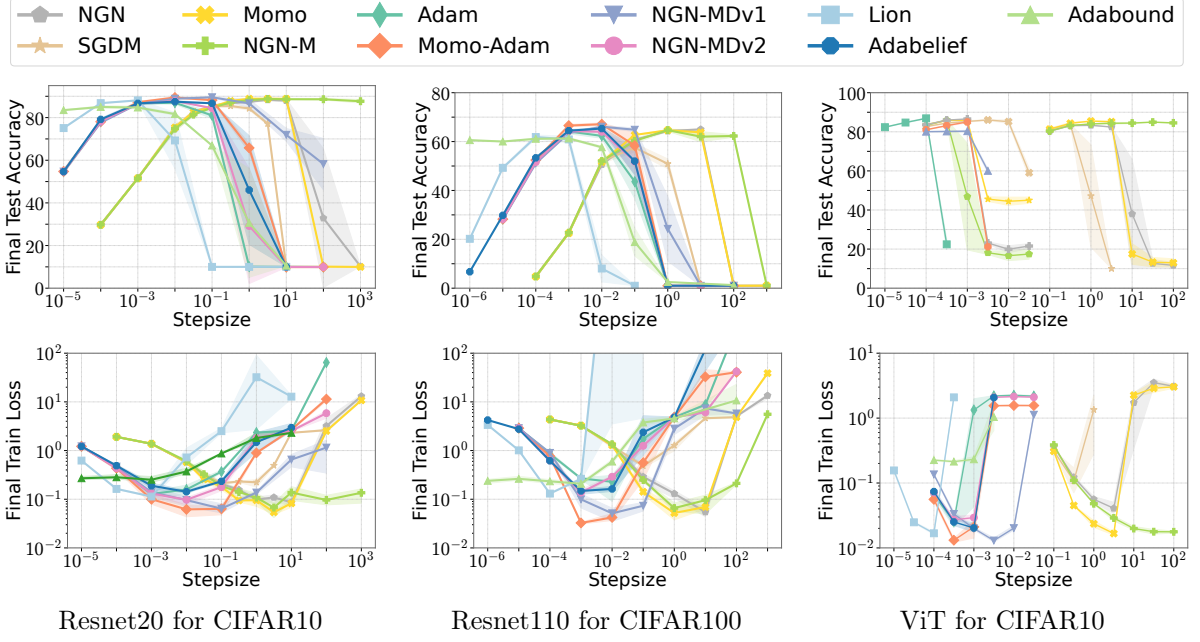


Figure J.5: Stability performance of various optimizers for (Resnet20, CIFAR10), (Resnet110, CIFAR100), (ViT, CIFAR10) workloads.

Table 9: Test perplexity of NGN-MDv1 when varying the learning rate and β_2 hyperparameters when training 70M language model on the FineWeb dataset.

lr	$\beta_2 = 0.6$	$\beta_2 = 0.8$	$\beta_2 = 0.9$	$\beta_2 = 0.95$	$\beta_2 = 0.999$
3e-4	51.8 \pm 0.6	49.2 \pm 0.4	47.8 \pm 0.3	47.0 \pm 0.2	47.0 \pm 0.2
1e-3	42.6 \pm 0.3	40.5 \pm 0.1	39.3 \pm 0.2	38.6 \pm 0.1	38.9 \pm 0.1
3e-3	39.4 \pm 0.2	37.5 \pm 0.2	36.3 \pm 0.1	35.9 \pm 0.1	36.5 \pm 0.4
1e-2	37.8 \pm 0.2	35.9 \pm 0.1	35.1 \pm 0.3	34.7 \pm 0.3	35.0 \pm 0.3
3e-2	37.8 \pm 0.3	35.8 \pm 0.1	34.9 \pm 0.1	34.6 \pm 0.1	250 \pm 50

Table 10: Test perplexity of Adam when varying the learning rate and β_2 hyperparameters when training 70M language model on the FineWeb dataset.

lr	$\beta_2 = 0.6$	$\beta_2 = 0.8$	$\beta_2 = 0.9$	$\beta_2 = 0.95$	$\beta_2 = 0.999$
3e-4	46.1 \pm 0.2	46.6 \pm 0.1	46.5 \pm 0.2	46.2 \pm 0.3	46.5 \pm 0.1
1e-3	38.8 \pm 0.1	39.0 \pm 0.2	38.9 \pm 0.1	38.5 \pm 0.1	39.5 \pm 0.6
3e-3	38.8 \pm 0.3	36.3 \pm 0.1	36.1 \pm 0.2	36.0 \pm 0.1	36.7 \pm 0.8
1e-2	35.4 \pm 0.2	35.0 \pm 0.1	34.9 \pm 0.3	36 \pm 2	41 \pm 3
3e-2	550 \pm 250	120 \pm 80	160 \pm 5	210 \pm 60	500 \pm 20

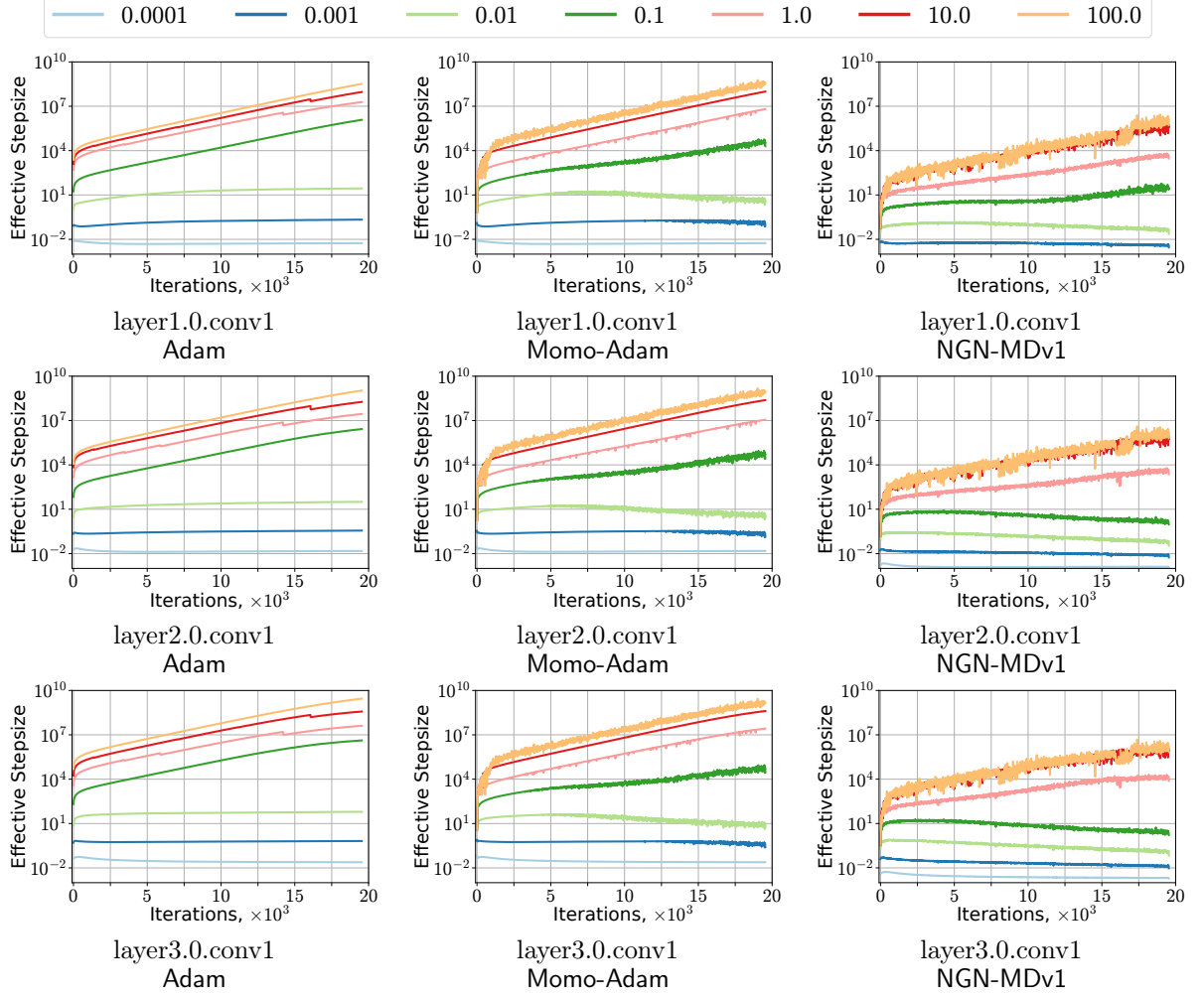


Figure J.6: The adaptive stepsize of Adam (**first column**), Momo-Adam (**second column**), and NGN-MDv1 (**third column**) algorithms in training ResNet20 model on CIFAR10 dataset. We plot the average stepsize $\frac{\gamma}{(\mathbf{D}_k)_{(j)}}$ (for Adam), $\frac{\tau_k}{(\mathbf{D}_k)_{(j)}}$ (for Momo-Adam), and $\frac{\gamma_k}{(\mathbf{D}_k)_{(j)}}$ (for NGN-MDv1) for the first convolution layer within each of 3 base blocks of ResNet20 architecture varying the stepsize hyperparameter of the algorithms (c for NGN-M and NGN, α_0 for Momo, and learning rate parameter for Adam).

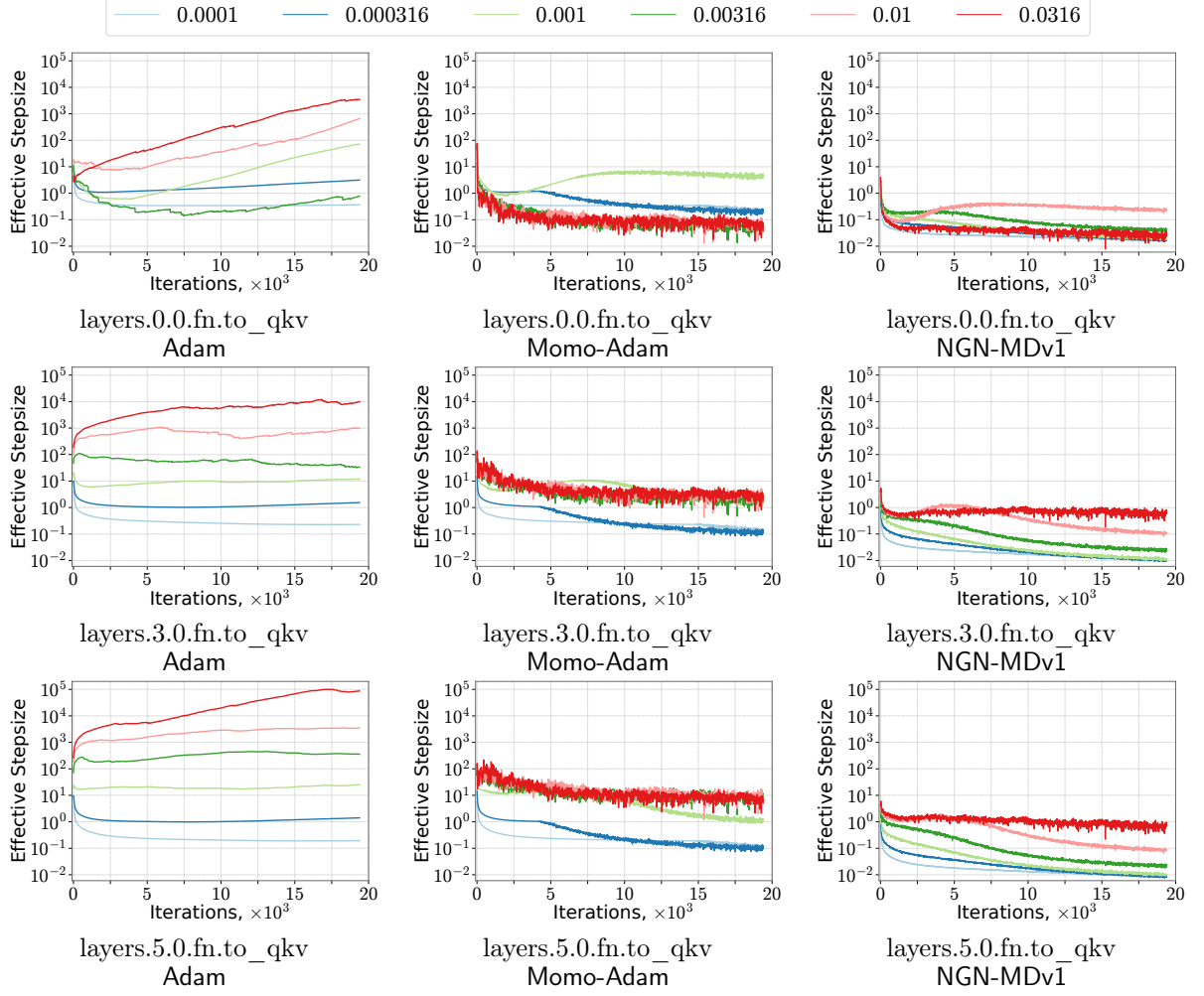


Figure J.7: The adaptive stepsize of Adam (**first column**), Momo-Adam (**second column**), and NGN-MDv1 (**third column**) algorithms in training ViT model on CIFAR10 dataset. We plot the average stepsize $\frac{\gamma}{(\mathbf{D}_k)_{(j)}}$ (for Adam), $\frac{\tau_k}{(\mathbf{D}_k)_{(j)}}$ (for Momo-Adam), and $\frac{\gamma_k}{(\mathbf{D}_k)_{(j)}}$ (for NGN-MDv1) for the attention layer within each of the first, fourth, and sixth base blocks of ViT architecture varying the step-size hyperparameter of the algorithms (c for NGN-M and NGN, α_0 for Momo, and learning rate parameter for Adam).

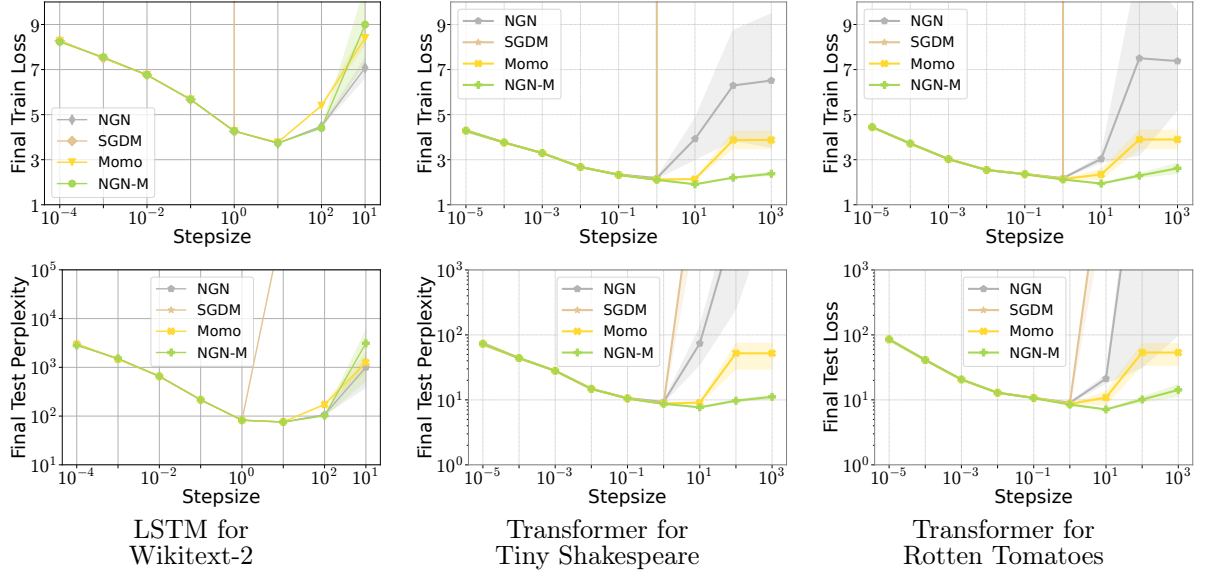


Figure J.8: Stability performance of algorithms supporting momentum and diagonal step-size varying step-size hyperparameter (c for NGN-M and NGN, α_0 for Momo, and step-size for SGDM). We observe that NGN-M achieves the training loss close to the best possible for a wider range of the step-size hyperparameter.

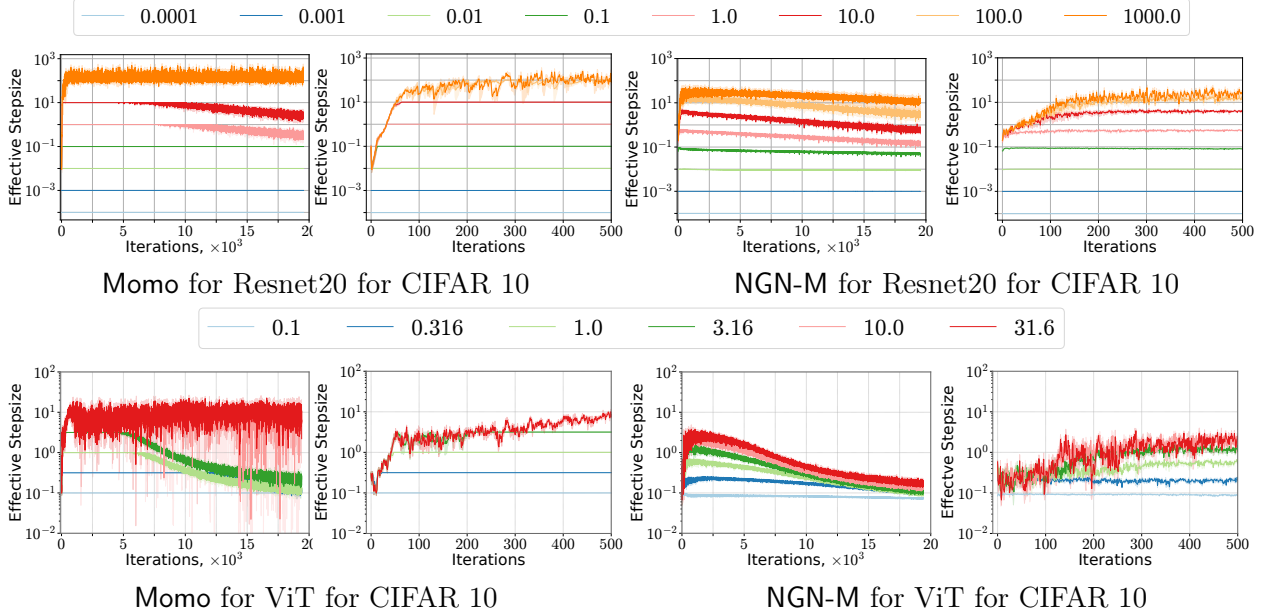


Figure J.9: The step-size of Momo and NGN-M during the training. We demonstrate the step-sizes τ_k for Momo and γ_k for NGN-M varying step-size parameters α_0 for Momo and c for NGN-M.

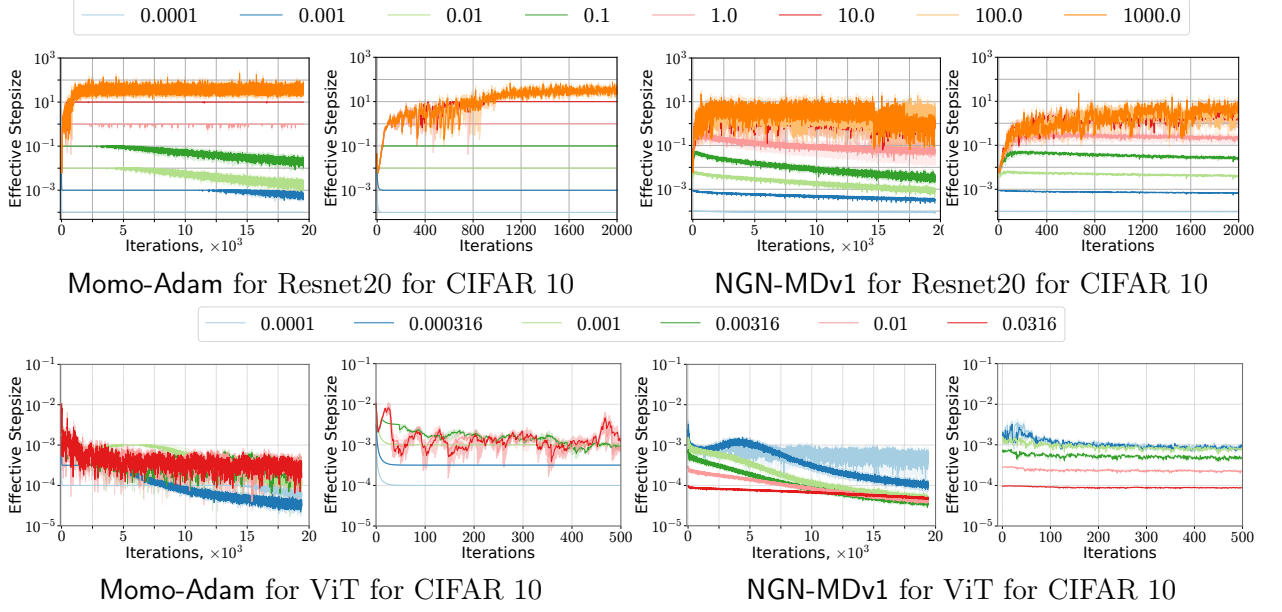


Figure J.10: The step-size of Momo-Adam and NGN-MDv1 during the training. We demonstrate the step-sizes τ_k for Momo-Adam and γ_k for NGN-MDv1 varying step-size parameters α_0 for Momo and c for NGN-MDv1.

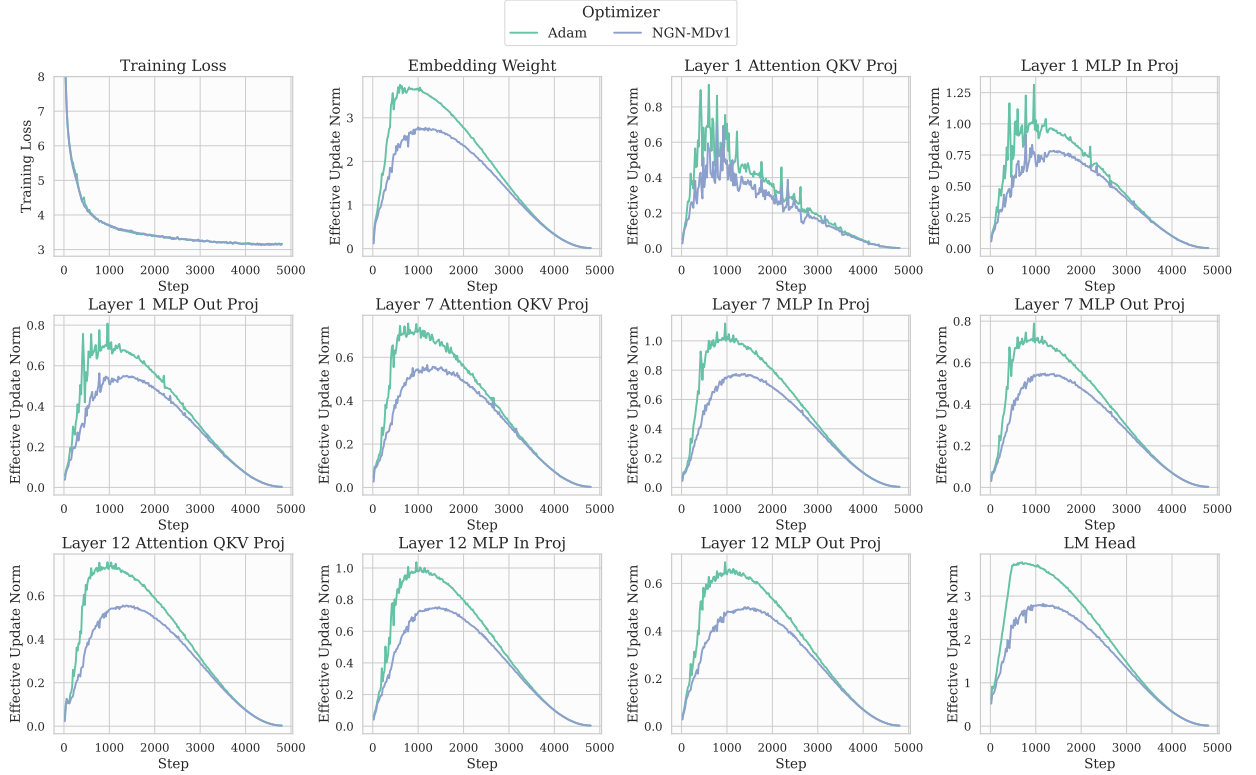


Figure J.11: Magnitude of updates when training 160M language model with Adam and NGN-MDv1 and step-size hyperparameter 0.003.

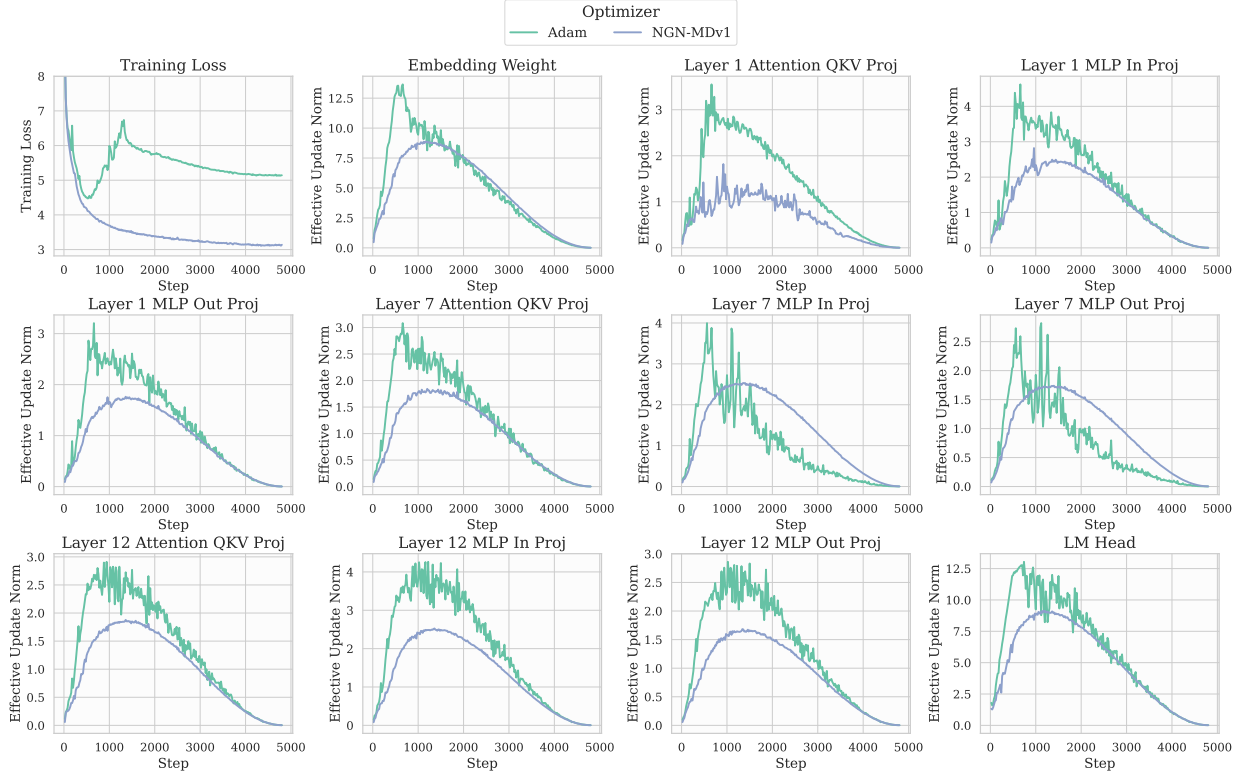


Figure J.12: Magnitude of updates when training 160M language model with Adam and NGN-MDv1 and step-size hyperparameter 0.01.

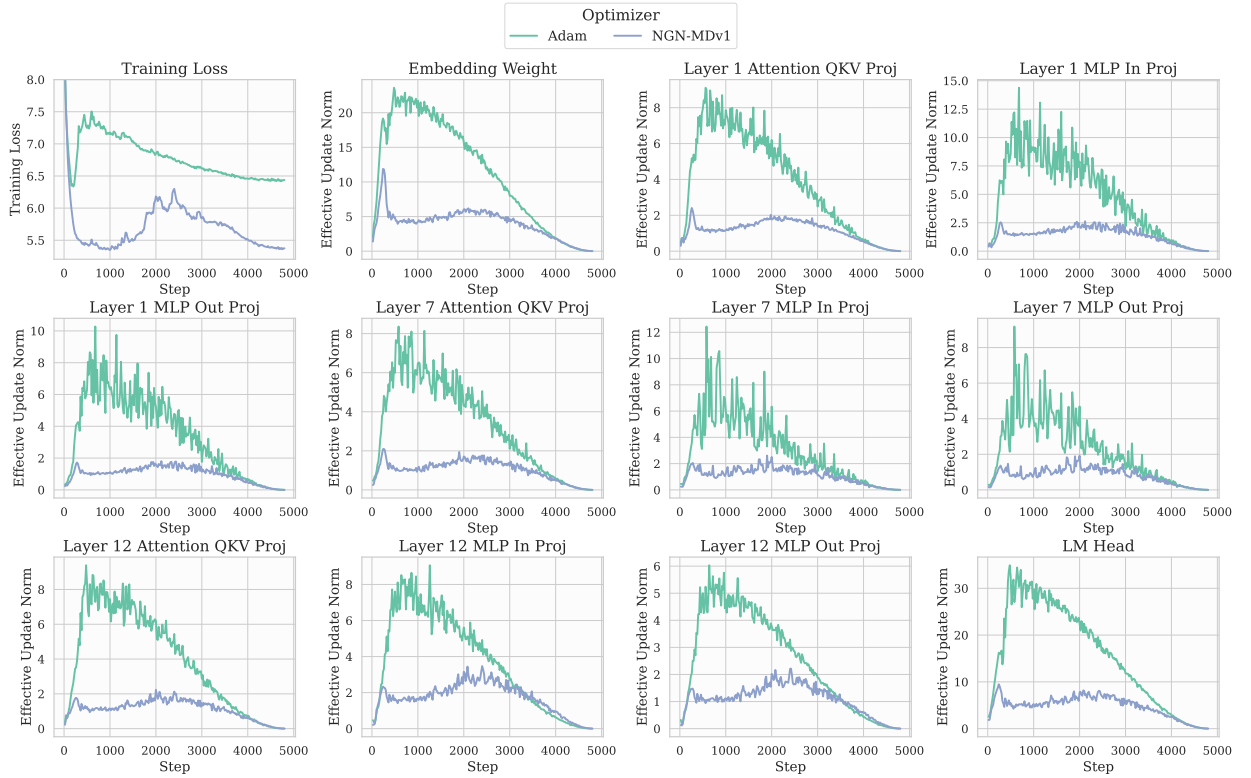


Figure J.13: Magnitude of updates when training 160M language model with Adam and NGN-MDv1 and step-size hyperparameter 0.03.

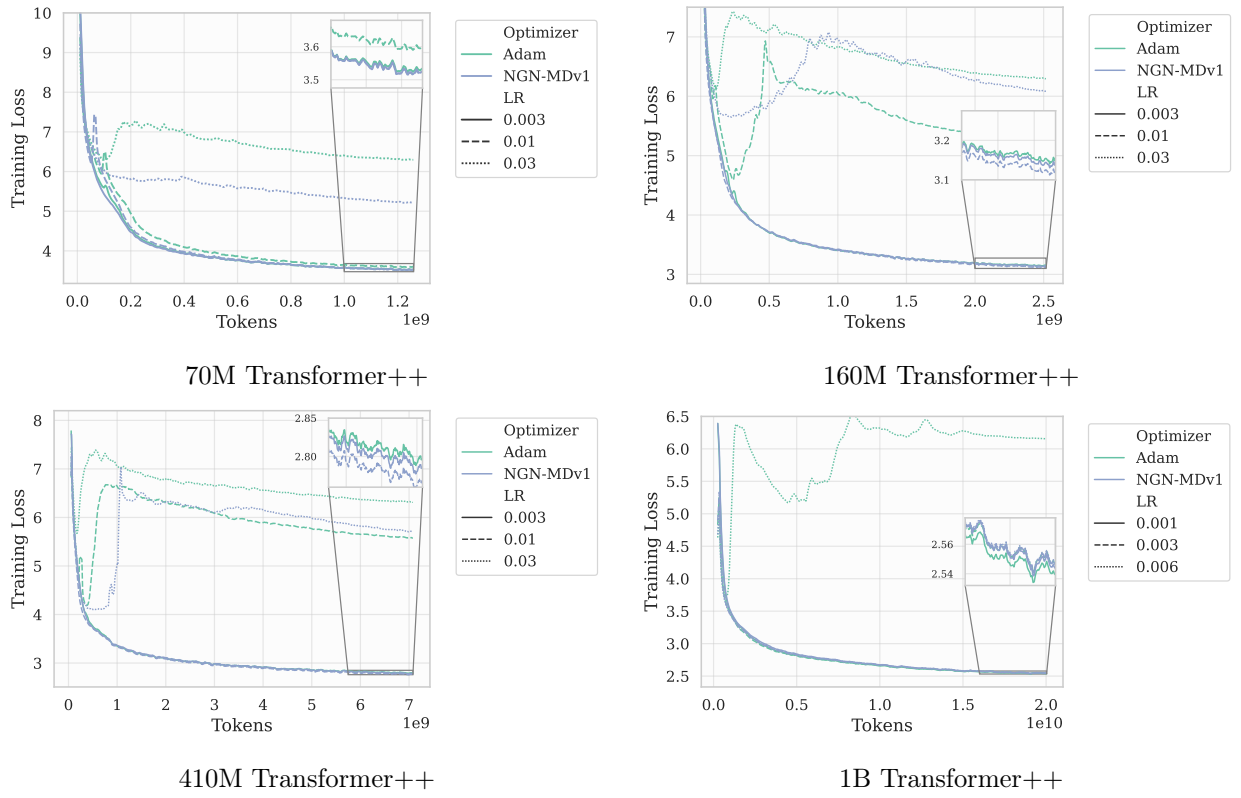


Figure J.14: Training dynamics when training language model at different sizes.