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

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Modern optimization algorithms that incorporate momentum and adaptive step-size
2 offer improved performance in numerous challenging deep learning tasks. However,
3 their effectiveness is often highly sensitive to the choice of hyperparameters, espe-
4 cially the step-size. Tuning these parameters is often difficult, resource-intensive,
5 and time-consuming. Therefore, recent efforts have been directed toward enhancing
6 the stability of optimizers across a wide range of hyperparameter choices [66]. In
7 this paper, we introduce an algorithm that matches the performance of state-of-the-
8 art optimizers while improving stability to the choice of the step-size hyperparam-
9 eter through a novel adaptation of the NGN step-size method [55]. Specifically, we
10 propose a momentum-based version (NGN-M) that attains the standard convergence
11 rate of $\mathcal{O}(1/\sqrt{K})$ under less restrictive assumptions, without the need for inter-
12 polation condition or assumptions of bounded stochastic gradients or iterates, in
13 contrast to previous approaches. Additionally, we empirically demonstrate that the
14 combination of the NGN step-size with momentum results in enhanced robustness
15 to the choice of the step-size hyperparameter while delivering performance that is
16 comparable to or surpasses other state-of-the-art optimizers.

17 1 Introduction

18 Adaptive methods such as Adam [36] and RMSprop [25] are widely used in machine learning
19 due to their established advantages over (momentum) SGD, particularly in tasks such as training
20 Transformers [7, 74, 75]. These methods adaptively scale the step-size across different dimensions
21 (parameters) based on their respective statistics, effectively acting as a diagonal preconditioning.

22 Although these methods perform well in practice, existing theoretical analyses typically require
23 stringent assumptions on the noise structure of the stochastic gradients, such as sub-Gaussian
24 noise [41] or affine noise models [76, 91]: Relaxing these assumptions remains an open challenge.
25 Another well-known issue of Adam is its performance sensitivity to the step-size hyperparameter [82,
26 9], particularly when training Transformers, where loss spikes are commonly observed [51, 83].
27 This often necessitates careful adjustments of the hyperparameters throughout the training process
28 [92, 10], which can be costly in terms of computational resources [54]. Consequently, there has
29 been growing interest in developing optimization methods that are more robust to hyperparameter
30 selection [66]. In addition to adapting the step-size, Adam and other state-of-the-art optimizers also
31 rely on momentum [59], a broadly used technique that has been shown to enhance performance both
32 theoretically [12, 17, 28] and practically [9, 19, 30]. Besides speeding up convergence, momentum is
33 known as a technique to reduce the variance of stochastic algorithms [48, 13], improving stability as
34 well as generalization in some settings [30].

35 In this work, we address the aforementioned drawbacks of Adam by developing a new algorithm
36 based on the recently proposed NGN step-size [55], an improved variant of the Stochastic Polyak Step-

size [45] 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 [55] 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 [22]. This issue is at the root of recent research focusing on transferring hyperparameters across architectures at different scales, therefore avoiding expensive tuning pipelines [85, 86, 6]. Yet, in the convex setting, choosing the learning rate can already be difficult – an issue that was studied already in Polyak [60] 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 [45, 56, 31], 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 [55] 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 [55] 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 [45] proves convergence to the exact solution only if the interpolation condition holds¹.

Polyak Step-size and Heavy-ball Momentum. Heavy-ball momentum methods, stemming from the work of Polyak [59], have gained significant attention over the years due to their benefits, including acceleration on convex quadratics [29, 40, 5], convex-like [78], and non-convex problems [12], as well as their variance reduction abilities [48, 13]. This has led to growing interest in the combination of Polyak step-size and heavy-ball momentum, which is an active area of research [2, 64, 2, 79, 53]. Recently, Schaipp et al. [66] 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 [52, 93]. Algorithms employing diagonal step-sizes, such as Adam and SignSGD [3], typically outperform non-diagonal methods in language modeling tasks by addressing issues such as class imbalance (where certain words appear more frequently than others) [38, 39] and heavy-tailed noise [89, 90, 11].

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 (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} [45]	$\mathcal{O}(1/K + \sigma_{\text{int}}^2)$	✗	✗	Conv. to non-vanishing neighbourhood
ALR-SMAG [79]	$\mathcal{O}((1 - \rho)^K + \sigma_{\text{int}}^2)$	✓	✗	Strong convexity Conv. to non-vanishing neighbourhood
Momo [66]	$\mathcal{O}(1/\sqrt{K})$	✓	✗	Bounded stoch. gradients Interpolation
Momo-Adam [66]	✗	✓	✓	Momo framework for Adam
MomSPS _{max} [53]	$\mathcal{O}(1/K + \sigma_{\text{int}}^2)$	✓	✗	Conv. to non-vanishing neighbourhood
NGN [55]	$\mathcal{O}(1/\sqrt{K})$	✗	✗	—
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]	$\mathcal{O}(1/\sqrt{K})$	✗	✓	—

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 [45, 79, 23, 53, 55]. Only a few works propose a diagonal-wise modification of Polyak-step-size by either using Adam preconditioner [66] as a weight matrix or incorporating second-order information from the objective function [43, 63].

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 [55], 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

$$p^k := \operatorname{argmin}_p \left[(r(x^k) + \nabla r(x^k)^\top p)^2 + \frac{1}{2c} \|p\|^2 \right]. \quad (1)$$

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 [55], 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 [55]). 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

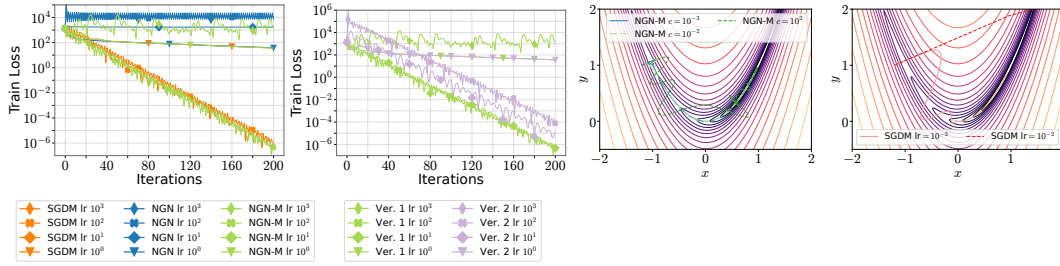


Figure 1: **Left:** Comparison of SGDM, NGN, NGN-M for linear regression on normalized Diabetes dataset varying a step-size hyperparameter. **Second left:** Comparison of two options on how momentum can be used in combination with NGN step-size. **Third and fourth:** Comparison of SGDM and NGN-M on the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

heavy-ball update [53]. 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 [79, 66]. 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 [71] 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 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 1.3 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

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

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  $S_k \subseteq [n]$ 
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

```

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

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

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)

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

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

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.

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 (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 [45, 56, 55, 53, 66].

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  $S_k \subseteq [n]$ 
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

```

175 **Assumption 4.1.** Each f_i is convex and L -smooth, i.e., for all $x, y \in \mathbb{R}^d$ and $i \in [n]$ we have
176 $\langle \nabla f_i(x), y - x \rangle \geq f_i(x) - f_i(y)$ and $\|\nabla f_i(x) - \nabla f_i(y)\| \leq L\|x - y\|$.

177 **Assumption 4.2.** 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
178 bounded. We say that the interpolation holds if $\sigma_{\text{int}}^2 = 0$.

179 4.2 Convergence Guarantees

180 **Theorem 4.3.** Let Assumptions 4.1, 4.2 hold. Let the step-size hyperparameter $c > 0$ and the
181 momentum parameter $\beta = \frac{\lambda}{1+\lambda}$ be constants where $\lambda \leq \min\{cL, 0.5(1+cL)^{-1}(1+2cL)^{-1}\}$. Then
182 the iterates of NGN-M (Alg. 1) satisfy

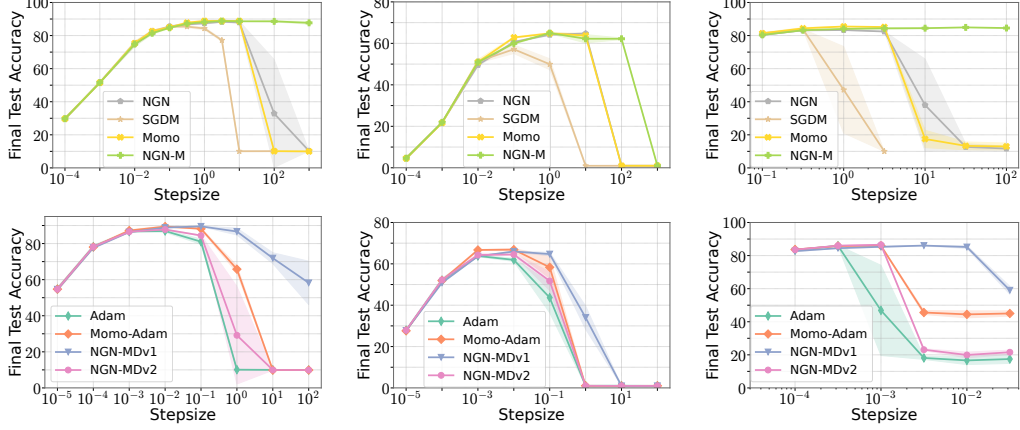
$$\mathbb{E}[f(\bar{x}^{K-1}) - f(x^*)] \leq \frac{\|x^0 - x^*\|^2 (1+2cL)^2}{cK} + 8cL(1+2cL)^2 \sigma_{\text{int}}^2 + 2cL \max\{2cL - 1, 0\} \sigma_{\text{pos}}^2,$$

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

185 Importantly, we show that (i) when the constant c is sufficiently small, NGN-M attains the same
186 convergence rate as SGDM [20]. Moreover, for any choice of c , we demonstrate that the NGN-M
187 iterates provably converge to a neighborhood of the optimum and thereafter remain within it; (ii)
188 Unlike prior works, our analysis does not rely on strong assumptions such as bounded gradients,
189 interpolation, or a bounded domain; (iii) For small values of c , NGN-M converges to the exact
190 solution while algorithms such as MomSPS and ALR-SMAG were shown to converge up to a non-
191 vanishing neighborhood of the solution only³. Notably, the non-vanishing neighborhood disappears
192 when the problem satisfies interpolation: We refer to Table 1 for more details and exact rates;
193 (iv) The momentum parameter β is theoretically recommended to be set sufficiently small. A
194 default value of $\beta = 0.9$ is commonly used and works well in our experiments. This discrepancy
195 between theoretical guidance and practical implementation has also been observed in prior works
196 on momentum [21, 44, 79, 78, 53]. Interestingly, for simple functions we can establish convergence
197 even when β is large (see Appendix F), indicating that the small- β requirement may be an artifact of
198 the existing proving techniques rather than an inherent algorithmic limitation of NGN-M. We leave
199 a comprehensive study of arbitrary β values across general convex objectives for future work; (v)
200 While Theorem 4.3 requires knowing the total iteration count K to ensure convergence, this might be
201 impractical: We therefore also prove convergence using a diminishing step-size of order $1/\sqrt{k}$ in
202 Appendix E, which removes the need to preset K ; (vi) Finally, we corroborate our analysis as we run
203 NGN-M with the theory-derived values of c to a quadratic problem that satisfies all our assumptions:
204 We observe NGN-M's rapid convergence with theoretical step-size hyperparameters in practice—see
205 Appendix I.3 and Figure I.4 therein.

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

³In fact, this is an inherited property of SPS analysis from [45].

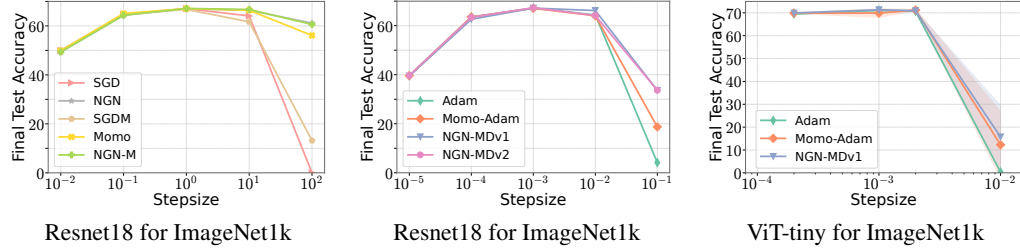


Resnet20 for CIFAR 10

Resnet10 for CIFAR 100

ViT for CIFAR 10

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.



Resnet18 for ImageNet1k

Resnet18 for ImageNet1k

ViT-tiny for ImageNet1k

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.

208 First, we make use of the Iterative Moving Average (IMA) formulation of momentum [67]. Specifi-
 209 cally, we define a sequence of virtual iterates $\{z^k\}$ whose update rule is of the form

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

210 Next, one of the key technical strategies we follow is splitting the step-size γ_k into two parts: a fixed
 211 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
 212 of the step-size γ_k enables us to regulate the balance between the descent term, which drives
 213 improvement in the objective, and the error term, which reflects possible inaccuracies. More precisely,
 214 the descent term is weighted by c while the error term proportional to σ_{int}^2 is weighted by c^2 , which
 215 suggests that c has to be chosen to tradeoff the two terms to lead to the exact convergence similarly to
 216 the standard analysis of SGD [20]. In contrast, MomSPS and Momo algorithms achieve the exact
 217 convergence only under the interpolation regime.

218 5 Experiments

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

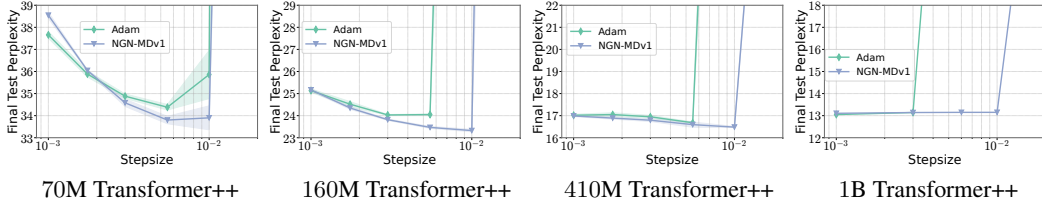


Figure 4: Language Modeling on SlimPajama. Stability comparison with respect to the step-size hyperparameter across different model sizes and optimizers. At all model capacities, NGN-MDv1 achieves similar or lower perplexity, showing better stability and improved performance at larger learning rates. 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.

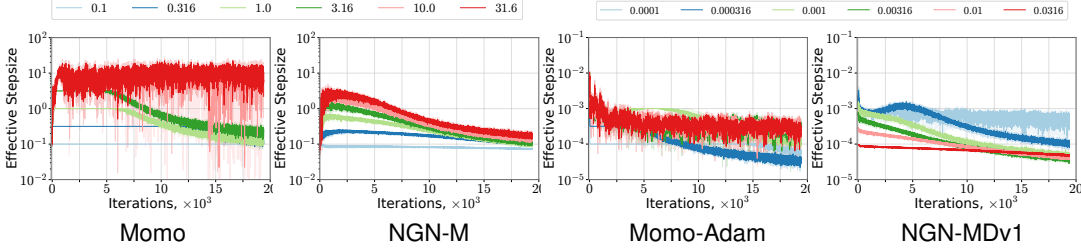


Figure 5: The step-size of Momo, NGN-M (two left), Momo-Adam and NGN-MDv1 (two right) during the training of ViT on CIFAR10. We demonstrate the step-sizes τ_k for Momo and Momo-Adam and γ_k for NGN-M and NGN-MDv1 varying step-size parameters α_0 and c correspondingly. We refer to Figures J.9 and J.10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

225 **Comparison on Standard Benchmarks.** First, we test the performance of NGN-M against other
 226 methods that use momentum, such as SGDM, Momo, MomSPS, ALR-SMAG, and NGN. The tests
 227 include the training of Resnet20 [24] and ViT [15] on the CIFAR10 dataset [37], and Resnet110
 228 on CIFAR100. Second, we test the performance of NGN-MD against Adam and Momo-Adam that –
 229 contrary to NGN-M – both use component-wise preconditioning. All experiments in this section do
 230 not use learning rate schedulers or weight decay.

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

241 **Vision Experiments on ImageNet.** Having observed promising results on workloads of small
 242 and medium size, we switch to larger tasks and datasets. We first train a ResNet18 on ImageNet1k
 243 [14]. This represents the first task in which we pair our proposed algorithms with a learning rate
 244 schedule. As illustrated in Figure 3 and Table 4, NGN-M achieves the highest validation accuracy,
 245 while exhibiting higher robustness across larger step-sizes, improving over both NGN and Momo.
 246 Among adaptive methods, NGN-MDv1 compares favorably against Adam and MomoAdam, while
 247 once again achieving higher performance on a wider range of learning rates (Table 5). Appendix J.4
 248 reports additional ablations on ImageNet32 and train loss stability results.

249 Finally, we test the effectiveness of the proposed algorithms on vision transformers [15]. These
 250 models are trained for a longer horizon compared to convolutional architectures, are notoriously
 251 sensitive to initial learning rate, and require adaptive step-sizes. We follow the protocol of Schaipp
 252 et al. [66], which includes cosine annealing, but without any weight decay regularization. As
 253 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 [61] with 70M, 160M, 410M, and 1B parameters around Chinchilla optimum [27] on SlimPajama-627B [72]. 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 in 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 [69] or ZeRO-2 [62]—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 [32] can be further improved by incorporating the NGN step-size.

References

- [1] 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 41 and 44)
- [2] 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 2)
- [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 2 and 23)
- [4] 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, Shivanshu 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 47)
- [5] 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 2)
- [6] 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)
- [7] Tom B Brown. Language models are few-shot learners. *arXiv preprint arXiv:2005.14165*, 2020. (Cited on page 1)
- [8] 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 50)
- [9] 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 page 1)
- [10] 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)
- [11] 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 2)
- [12] Ashok Cutkosky and Harsh Mehta. Momentum improves normalized sgd. In *International conference on machine learning*. PMLR, 2020. (Cited on pages 1 and 2)
- [13] Ashok Cutkosky and Francesco Orabona. Momentum-based variance reduction in non-convex sgd. *Advances in neural information processing systems*, 2019. (Cited on pages 1 and 2)
- [14] 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 8)
- [15] 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 page 8)

- [16] 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 53)
- [17] Ilyas Fatkhullin, Alexander Tyurin, and Peter Richtárik. Momentum provably improves error feedback! *Advances in Neural Information Processing Systems*, 2024. (Cited on page 1)
- [18] 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 38)
- [19] 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 1)
- [20] Guillaume Garrigos and Robert M Gower. Handbook of convergence theorems for (stochastic) gradient methods. *arXiv preprint arXiv:2301.11235*, 2023. (Cited on pages 6, 7, 23, 25, 28, and 32)
- [21] 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 6)
- [22] Ian J. Goodfellow, Yoshua Bengio, and Aaron Courville. *Deep Learning*. MIT Press, 2016. (Cited on page 2)
- [23] 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)
- [24] 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)
- [25] 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)
- [26] Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. *Neural Computation*, 1997. (Cited on page 49)
- [27] 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 9)
- [28] Rustem Islamov, Yuan Gao, and Sebastian U Stich. Near optimal decentralized optimization with compression and momentum tracking. *arXiv preprint arXiv:2405.2011*, 2024. (Cited on page 1)
- [29] 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 2)
- [30] Samy Jelassi and Yuanzhi Li. Towards understanding how momentum improves generalization in deep learning. In *International Conference on Machine Learning*, 2022. (Cited on page 1)
- [31] 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)
- [32] 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 9)

- [33] 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 48)
- [34] Andrej Karpathy. char-rnn. <https://github.com/karpathy/char-rnn>, 2015. (Cited on page 49)
- [35] Andrej Karpathy. Nanogpt. <https://github.com/karpathy/nanoGPT>, 2022. (Cited on pages 47 and 49)
- [36] Diederik Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *International Conference on Learning Representations (ICLR)*, 2015. (Cited on pages 1 and 53)
- [37] Alex Krizhevsky, Vinod Nair, and Geoffrey Hinton. Learning multiple layers of features from tiny images. *Scientific Report*, 2009. (Cited on page 8)
- [38] 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 2)
- [39] 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 2)
- [40] 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 2)
- [41] Haochuan Li, Alexander Rakhlin, and Ali Jadbabaie. Convergence of adam under relaxed assumptions. *Advances in Neural Information Processing Systems*, 2024. (Cited on page 1)
- [42] 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 41)
- [43] 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)
- [44] 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 6)
- [45] 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, 5, and 6)
- [46] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. *arXiv preprint arXiv: 1711.05101*, 2019. (Cited on page 42)
- [47] 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 50)
- [48] Jerry Ma and Denis Yarats. Quasi-hyperbolic momentum and adam for deep learning. *arXiv preprint arXiv:1810.06801*, 2018. (Cited on pages 1 and 2)
- [49] 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 49)
- [50] 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 49)

- [51] 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)
- [52] 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 2)
- [53] 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 2, 3, 4, 5, 6, and 23)
- [54] 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)
- [55] 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, 5, 23, and 24)
- [56] 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 5)
- [57] 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 49)
- [58] 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 47)
- [59] Boris T. Polyak. Some methods of speeding up the convergence of iteration methods. *Ussr computational mathematics and mathematical physics*, 1964. (Cited on pages 1, 2, and 38)
- [60] Boris T Polyak. Introduction to optimization. *New York, Optimization Software*, 1987. (Cited on page 2)
- [61] Alec Radford, Karthik Narasimhan, Tim Salimans, and Ilya Sutskever. Improving language understanding by generative pre-training. Technical report, OpenAI, 2019. (Cited on page 9)
- [62] 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 9 and 41)
- [63] 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)
- [64] Samer Saab, Shashi Phoha, Minghui Zhu, and Asok Ray. An adaptive polyak heavy-ball method. *Machine Learning*, 2022. (Cited on page 2)
- [65] Mher Safaryan and Peter Richtárik. Stochastic sign descent methods: New algorithms and better theory. In *International Conference on Machine Learning*, 2021. (Cited on page 23)
- [66] 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, 2, 3, 4, 5, 8, and 47)
- [67] 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 23)
- [68] Noam Shazeer. Glu variants improve transformer. *arXiv preprint arXiv: 2002.05202*, 2020. (Cited on page 47)

- [69] 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 9 and 41)
- [70] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*, 2014. (Cited on page 49)
- [71] 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)
- [72] 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 9 and 48)
- [73] Jianlin Su, Yu Lu, Shengfeng Pan, Ahmed Murtadha, Bo Wen, and Yunfeng Liu. Roformer: Enhanced transformer with rotary position embedding, 2023. (Cited on page 47)
- [74] 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)
- [75] 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 47)
- [76] 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)
- [77] 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 38)
- [78] 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 2 and 6)
- [79] 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 2, 3, 4, and 6)
- [80] Rachel Ward, Xiaoxia Wu, and Leon Bottou. Adagrad stepsizes: Sharp convergence over nonconvex landscapes. *Journal of Machine Learning Research*, 2020. (Cited on page 24)
- [81] Ross Wightman. Pytorch image models. <https://github.com/rwightman/pytorch-image-models>, 2019. (Cited on page 47)
- [82] 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)
- [83] 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)
- [84] Lechao Xiao. Rethinking conventional wisdom in machine learning: From generalization to scaling. *arXiv preprint arXiv: 2409.15156*, 2024. (Cited on pages 41 and 44)

- [85] 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)
- [86] Greg Yang, James B Simon, and Jeremy Bernstein. A spectral condition for feature learning. *arXiv preprint arXiv:2310.17813*, 2023. (Cited on page 2)
- [87] Biao Zhang and Rico Sennrich. Root mean square layer normalization, 2019. (Cited on page 47)
- [88] Guodong Zhang, Chaoqi Wang, Bowen Xu, and Roger Grosse. Three mechanisms of weight decay regularization. *arXiv preprint arXiv:1810.12281*, 2018. (Cited on page 41)
- [89] 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 2)
- [90] 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 2)
- [91] 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*, 2024. (Cited on page 1)
- [92] 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)
- [93] 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*, 2024. (Cited on page 2)
- [94] 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 50)

NeurIPS Paper Checklist

1. Claims

Question: Do the main claims made in the abstract and introduction accurately reflect the paper’s contributions and scope?

Answer: [Yes]

Justification: We propose a novel scheme how to combine NGN step-size with momentum and component-wise step-size. We provide extensive theoretical and numerical analysis of NGN-M and NGN-MD to support the claims made in the abstract and introduction.

Guidelines:

- The answer NA means that the abstract and introduction do not include the claims made in the paper.
- The abstract and/or introduction should clearly state the claims made, including the contributions made in the paper and important assumptions and limitations. A No or NA answer to this question will not be perceived well by the reviewers.
- The claims made should match theoretical and experimental results, and reflect how much the results can be expected to generalize to other settings.
- It is fine to include aspirational goals as motivation as long as it is clear that these goals are not attained by the paper.

2. Limitations

588 Question: Does the paper discuss the limitations of the work performed by the authors?

589 Answer: [\[Yes\]](#)

590 Justification: We include a discussion on the limitations of our work in the conclusion and
591 appendix.

592 Guidelines:

- 593 • The answer NA means that the paper has no limitation while the answer No means that
594 the paper has limitations, but those are not discussed in the paper.
- 595 • The authors are encouraged to create a separate "Limitations" section in their paper.
- 596 • The paper should point out any strong assumptions and how robust the results are to
597 violations of these assumptions (e.g., independence assumptions, noiseless settings,
598 model well-specification, asymptotic approximations only holding locally). The authors
599 should reflect on how these assumptions might be violated in practice and what the
600 implications would be.
- 601 • The authors should reflect on the scope of the claims made, e.g., if the approach was
602 only tested on a few datasets or with a few runs. In general, empirical results often
603 depend on implicit assumptions, which should be articulated.
- 604 • The authors should reflect on the factors that influence the performance of the approach.
605 For example, a facial recognition algorithm may perform poorly when image resolution
606 is low or images are taken in low lighting. Or a speech-to-text system might not be
607 used reliably to provide closed captions for online lectures because it fails to handle
608 technical jargon.
- 609 • The authors should discuss the computational efficiency of the proposed algorithms
610 and how they scale with dataset size.
- 611 • If applicable, the authors should discuss possible limitations of their approach to
612 address problems of privacy and fairness.
- 613 • While the authors might fear that complete honesty about limitations might be used by
614 reviewers as grounds for rejection, a worse outcome might be that reviewers discover
615 limitations that aren't acknowledged in the paper. The authors should use their best
616 judgment and recognize that individual actions in favor of transparency play an impor-
617 tant role in developing norms that preserve the integrity of the community. Reviewers
618 will be specifically instructed to not penalize honesty concerning limitations.

619 3. Theory assumptions and proofs

620 Question: For each theoretical result, does the paper provide the full set of assumptions and
621 a complete (and correct) proof?

622 Answer: [\[Yes\]](#)

623 Justification: We provide all assumptions used in the analysis in Section 4 and Appendix C.
624 The proofs of convergence and stability are deferred to Appendix C, D, E, and F.

625 Guidelines:

- 626 • The answer NA means that the paper does not include theoretical results.
- 627 • All the theorems, formulas, and proofs in the paper should be numbered and cross-
628 referenced.
- 629 • All assumptions should be clearly stated or referenced in the statement of any theorems.
- 630 • The proofs can either appear in the main paper or the supplemental material, but if
631 they appear in the supplemental material, the authors are encouraged to provide a short
632 proof sketch to provide intuition.
- 633 • Inversely, any informal proof provided in the core of the paper should be complemented
634 by formal proofs provided in appendix or supplemental material.
- 635 • Theorems and Lemmas that the proof relies upon should be properly referenced.

636 4. Experimental result reproducibility

637 Question: Does the paper fully disclose all the information needed to reproduce the main ex-
638 perimental results of the paper to the extent that it affects the main claims and/or conclusions
639 of the paper (regardless of whether the code and data are provided or not)?

640 Answer: [\[Yes\]](#)

Justification: We provide the Pytorch-based implementation in the supplementary. We use an open-source implementation of models and public datasets in our experiments. All training details are reported in Appendix J.

Guidelines:

- The answer NA means that the paper does not include experiments.
- If the paper includes experiments, a No answer to this question will not be perceived well by the reviewers: Making the paper reproducible is important, regardless of whether the code and data are provided or not.
- If the contribution is a dataset and/or model, the authors should describe the steps taken to make their results reproducible or verifiable.
- Depending on the contribution, reproducibility can be accomplished in various ways. For example, if the contribution is a novel architecture, describing the architecture fully might suffice, or if the contribution is a specific model and empirical evaluation, it may be necessary to either make it possible for others to replicate the model with the same dataset, or provide access to the model. In general, releasing code and data is often one good way to accomplish this, but reproducibility can also be provided via detailed instructions for how to replicate the results, access to a hosted model (e.g., in the case of a large language model), releasing of a model checkpoint, or other means that are appropriate to the research performed.
- While NeurIPS does not require releasing code, the conference does require all submissions to provide some reasonable avenue for reproducibility, which may depend on the nature of the contribution. For example
 - (a) If the contribution is primarily a new algorithm, the paper should make it clear how to reproduce that algorithm.
 - (b) If the contribution is primarily a new model architecture, the paper should describe the architecture clearly and fully.
 - (c) If the contribution is a new model (e.g., a large language model), then there should either be a way to access this model for reproducing the results or a way to reproduce the model (e.g., with an open-source dataset or instructions for how to construct the dataset).
 - (d) We recognize that reproducibility may be tricky in some cases, in which case authors are welcome to describe the particular way they provide for reproducibility. In the case of closed-source models, it may be that access to the model is limited in some way (e.g., to registered users), but it should be possible for other researchers to have some path to reproducing or verifying the results.

5. Open access to data and code

Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material?

Answer: [Yes]

Justification: We use only publicly available code and datasets and provide links to them in Appendix J.

Guidelines:

- The answer NA means that paper does not include experiments requiring code.
- Please see the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- While we encourage the release of code and data, we understand that this might not be possible, so “No” is an acceptable answer. Papers cannot be rejected simply for not including code, unless this is central to the contribution (e.g., for a new open-source benchmark).
- The instructions should contain the exact command and environment needed to run to reproduce the results. See the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- The authors should provide instructions on data access and preparation, including how to access the raw data, preprocessed data, intermediate data, and generated data, etc.

- The authors should provide scripts to reproduce all experimental results for the new proposed method and baselines. If only a subset of experiments are reproducible, they should state which ones are omitted from the script and why.
- At submission time, to preserve anonymity, the authors should release anonymized versions (if applicable).
- Providing as much information as possible in supplemental material (appended to the paper) is recommended, but including URLs to data and code is permitted.

6. Experimental setting/details

Question: Does the paper specify all the training and test details (e.g., data splits, hyper-parameters, how they were chosen, type of optimizer, etc.) necessary to understand the results?

Answer: [\[Yes\]](#)

Justification: We provide all training details in Appendix J and Table 3.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The experimental setting should be presented in the core of the paper to a level of detail that is necessary to appreciate the results and make sense of them.
- The full details can be provided either with the code, in appendix, or as supplemental material.

7. Experiment statistical significance

Question: Does the paper report error bars suitably and correctly defined or other appropriate information about the statistical significance of the experiments?

Answer: [\[Yes\]](#)

Justification: For each experiment, we plot the mean and a standard deviation across at least 3 random seeds (if the opposite is not stated). We do not include error bars when training 1B language model due to a limited resource availability which is necessary to run this experiment.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The authors should answer "Yes" if the results are accompanied by error bars, confidence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.
- The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, random drawing of some parameter, or overall run with given experimental conditions).
- The method for calculating the error bars should be explained (closed form formula, call to a library function, bootstrap, etc.)
- The assumptions made should be given (e.g., Normally distributed errors).
- It should be clear whether the error bar is the standard deviation or the standard error of the mean.
- It is OK to report 1-sigma error bars, but one should state it. The authors should preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis of Normality of errors is not verified.
- For asymmetric distributions, the authors should be careful not to show in tables or figures symmetric error bars that would yield results that are out of range (e.g. negative error rates).
- If error bars are reported in tables or plots, The authors should explain in the text how they were calculated and reference the corresponding figures or tables in the text.

8. Experiments compute resources

Question: For each experiment, does the paper provide sufficient information on the computer resources (type of compute workers, memory, time of execution) needed to reproduce the experiments?

Answer: [Yes]

Justification: We provide a description of compute resource we used for each experiment in Appendix J.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The paper should indicate the type of compute workers CPU or GPU, internal cluster, or cloud provider, including relevant memory and storage.
- The paper should provide the amount of compute required for each of the individual experimental runs as well as estimate the total compute.
- The paper should disclose whether the full research project required more compute than the experiments reported in the paper (e.g., preliminary or failed experiments that didn't make it into the paper).

9. Code of ethics

Question: Does the research conducted in the paper conform, in every respect, with the NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines>?

Answer: [Yes]

Justification: The research is conducted following the NeurIPS instructions, including the regulations regarding anonymity.

Guidelines:

- The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- If the authors answer No, they should explain the special circumstances that require a deviation from the Code of Ethics.
- The authors should make sure to preserve anonymity (e.g., if there is a special consideration due to laws or regulations in their jurisdiction).

10. Broader impacts

Question: Does the paper discuss both potential positive societal impacts and negative societal impacts of the work performed?

Answer: [Yes]

Justification: This paper presents work whose goal is to advance the field of Machine Learning. There are many potential societal consequences of our work, none of which we feel must be specifically highlighted here.

Guidelines:

- The answer NA means that there is no societal impact of the work performed.
- If the authors answer NA or No, they should explain why their work has no societal impact or why the paper does not address societal impact.
- Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.
- The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
- The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.

- If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

11. Safeguards

Question: Does the paper describe safeguards that have been put in place for responsible release of data or models that have a high risk for misuse (e.g., pretrained language models, image generators, or scraped datasets)?

Answer: [NA]

Justification:

Guidelines:

- The answer NA means that the paper poses no such risks.
- Released models that have a high risk for misuse or dual-use should be released with necessary safeguards to allow for controlled use of the model, for example by requiring that users adhere to usage guidelines or restrictions to access the model or implementing safety filters.
- Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
- We recognize that providing effective safeguards is challenging, and many papers do not require this, but we encourage authors to take this into account and make a best faith effort.

12. Licenses for existing assets

Question: Are the creators or original owners of assets (e.g., code, data, models), used in the paper, properly credited and are the license and terms of use explicitly mentioned and properly respected?

Answer: [NA]

Justification:

Guidelines:

- The answer NA means that the paper does not use existing assets.
- The authors should cite the original paper that produced the code package or dataset.
- The authors should state which version of the asset is used and, if possible, include a URL.
- The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.
- If assets are released, the license, copyright information, and terms of use in the package should be provided. For popular datasets, paperswithcode.com/datasets has curated licenses for some datasets. Their licensing guide can help determine the license of a dataset.
- For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.
- If this information is not available online, the authors are encouraged to reach out to the asset's creators.

13. New assets

Question: Are new assets introduced in the paper well documented and is the documentation provided alongside the assets?

Answer: [NA]

Justification:

Guidelines:

- The answer NA means that the paper does not release new assets.

- Researchers should communicate the details of the dataset/code/model as part of their submissions via structured templates. This includes details about training, license, limitations, etc.
- The paper should discuss whether and how consent was obtained from people whose asset is used.
- At submission time, remember to anonymize your assets (if applicable). You can either create an anonymized URL or include an anonymized zip file.

14. Crowdsourcing and research with human subjects

Question: For crowdsourcing experiments and research with human subjects, does the paper include the full text of instructions given to participants and screenshots, if applicable, as well as details about compensation (if any)?

Answer: [NA]

Justification:

Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Including this information in the supplemental material is fine, but if the main contribution of the paper involves human subjects, then as much detail as possible should be included in the main paper.
- According to the NeurIPS Code of Ethics, workers involved in data collection, curation, or other labor should be paid at least the minimum wage in the country of the data collector.

15. Institutional review board (IRB) approvals or equivalent for research with human subjects

Question: Does the paper describe potential risks incurred by study participants, whether such risks were disclosed to the subjects, and whether Institutional Review Board (IRB) approvals (or an equivalent approval/review based on the requirements of your country or institution) were obtained?

Answer: [NA]

Justification:

Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Depending on the country in which research is conducted, IRB approval (or equivalent) may be required for any human subjects research. If you obtained IRB approval, you should clearly state this in the paper.
- We recognize that the procedures for this may vary significantly between institutions and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the guidelines for their institution.
- For initial submissions, do not include any information that would break anonymity (if applicable), such as the institution conducting the review.

16. Declaration of LLM usage

Question: Does the paper describe the usage of LLMs if it is an important, original, or non-standard component of the core methods in this research? Note that if the LLM is used only for writing, editing, or formatting purposes and does not impact the core methodology, scientific rigor, or originality of the research, declaration is not required.

Answer: [NA]

Justification:

Guidelines:

- The answer NA means that the core method development in this research does not involve LLMs as any important, original, or non-standard components.
- Please refer to our LLM policy (<https://neurips.cc/Conferences/2025/LLM>) for what should or should not be described.

Appendix

901	Contents	
902	A Equivalent Formulations of NGN-M	23
903	B Technical Lemmas and Definitions	23
904	C Convergence of NGN-D	24
905	C.1 Convergence in General Non-convex Setting	25
906	C.2 Convergence under PL-condition	27
907	D Convergence of NGN-M	28
908	E Convergence of NGN-M with Decaying Step-size	32
909	F Stability of NGN-M on a Simple Problem	36
910	G How to Derive Diagonal NGN-based Step-size?	40
911	G.1 Design Comparison of NGN-MDv1 and NGN-MDv2	40
912	G.2 Computation Cost of NGN-MD	41
913	H How to add weight decay to NGN-MDv1?	41
914	H.1 Combining NGN-MDv1 and Weight Decay Regularization	42
915	H.2 Empirical Validation of the Proposed Combinations	44
916	I Additional Experiments on Toy Problems	45
917	I.1 Additional Experiments on the Problem with Many Minima	45
918	I.2 Comparison on Rosenbrock Function	45
919	I.3 Comparison on Quadratic Function with Theoretical Step-size	46
920	J Additional Experiments and Training Details	47
921	J.1 Training Details	47
922	J.2 Comparison Algorithms that Support Momentum	49
923	J.3 Comparison of Algorithms that Support Momentum and Diagonal Step-size	49
924	J.4 Additional ImageNet Experiments	49
925	J.5 Additional Comparison against Lion, Adabelief, Adabound	50
926	J.6 Comparison of Adaptive Step-sizes of Adam, Momo-Adam, and NGN-MDv1	50
927	J.7 Extended Comparison of Momentum-based Algorithms on NLP Tasks	53
928	J.8 Comparison of Algorithms with Diagonal Step-size	53
929	J.9 Effective Step-size of NGN-M, Momo, NGN-MDv1, and Momo-Adam	54
930	J.10 Effective Updates in Training Language Models	55

932 **A Equivalent Formulations of NGN-M**

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

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

934 We can rewrite the update rule using Iterative-Moving Average (IMA) approach presented in Proposi-
935 tion 1.6, Sebbouh et al. [67].

936 **Lemma A.1** (Proposition C.8 [53], Lemma 7.3 in [20]). *The iterates $\{x^k\}$ generated by NGN-M are*
937 *equivalent to the sequence $\{x^k\}$ generated by IMA update*

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

938 where

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

939 *Proof.* Let the sequences $\{x^k\}$ and $\{z^k\}$ be defined according to Equation (5). Let β be defined as
940 $\frac{\lambda}{1 + \lambda}$. Then we have

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

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

942 □

943 **B Technical Lemmas and Definitions**

944 **Definition B.1.** We say that the function ϕ admits \mathbf{L} -smooth with parameters $\mathbf{L} :=$
945 $(L_1, \dots, L_d), L_j \geq 0 \forall j \in [d]$, if the following inequality holds for all $x, h \in \mathbb{R}^d$

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

946 **Remark B.2.** If we set for all $j \in [d]$ $L_j := L$ then Definition B.1 reduces to standard L -smoothness.

947 This assumption is typically used in the context of coordinate adaptive algorithms such as SignSGD
948 [3, 65].

949 **Definition B.3.** The function $\phi: \mathbb{R}^d \rightarrow \mathbb{R}$ satisfies PL -condition with constant $\mu > 0$ if for all
950 $x, y \in \mathbb{R}^d$ we have

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

951 **Assumption B.4.** We assume that the coordinate-wise variance of the stochastic estimator is bounded,
952 i.e. for all $x \in \mathbb{R}^d$ and $j \in [d]$ we have

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

953 **Lemma B.5** (Lemma 4.9 from [55]). *Let each f_i be L -smooth for all i , then the step-size of NGN*
954 *satisfies*

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

973 **Theorem C.2.** Let f satisfies PL-condition and each f_i satisfies Definition B.1. Assume that
 974 Assumption B.4 holds. Then the iterates of NGN-D (Alg. 3) with step-size parameters $\{c_j\}_{j=1}^d$ such
 975 that $c_j \leq \min\{1/2L_j, 6/\mu\}$ satisfy

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

976 where $c_{\min} := \min_{j \in [d]} c_j$. Moreover, if $c_j = \mathcal{O}(\varepsilon)$ for all $j \in [d]$ then after $K =$
 977 $\max\{\mathcal{O}(\varepsilon^{-1}), \mathcal{O}(1)\} \log \varepsilon^{-1}$ iterations we obtain $\mathbb{E} [f(x^K) - f^*] \leq \mathcal{O}(\varepsilon)$.

978 To the best of our knowledge, this is the first result of the convergence of the Polyak-like step-size
 979 algorithm under the PL-condition. The convergence guarantees are similar to that of SGD [20].

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

981 **Lemma C.3** (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 .
 982 Let $f_{S_k}(x)$ satisfy Definition (B.1). Consider γ_j^k as in NGN-D (Algorithm 3), then we have

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

983 *Proof.* From Lemma B.7 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
 984 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,

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

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

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

986 and

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

987 that concludes the proof. \square

988 **Lemma C.4** (Fundamental Equality). Consider γ_j^k as in NGN-D (Algorithm 3). Then the following
 989 equality holds

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

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

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

991 which one can rewrite as

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

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

993 C.1 Convergence in General Non-convex Setting

994 **Theorem C.1.** Let each f_i satisfies Definition B.1. Assume that Assumption B.4 holds. Then the
 995 iterates of NGN-D (Alg. 3) with step-size parameters $\{c_j\}_{j=1}^d$ such that $c_j \leq 1/2L_j$ satisfy

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

996 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
 997 obtain $\min_{0 \leq k < K} \mathbb{E} [\|\nabla f(x^k)\|^2] \leq \mathcal{O}(\varepsilon^2)$.

998 *Proof.* First, we write separable Definition B.1

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

999 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
1000 conditional expectation with respect to x^k , as in this case we would have to analyze the product
1001 $\gamma_j^k \nabla_j f_{S_k}(x^k)$. Given bounds of the step-size γ_j^k from Lemma C.3, we can write the step-size as
1002 follows

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

1003 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
1004 of γ_j^k . Thus, we continue as follows

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

1005 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

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

1006 Therefore, we get

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

1007 We plug in (18) into (17) and get

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

1008 If $c_j \leq \frac{1}{2L_j}$, we get

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

1009

□

1010 We continue as follows

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

1011 Taking full expectation and unrolling the recursion above for all iterations $\{0, \dots, K-1\}$. Thus, we
1012 obtain

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

1013 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
1014 step-size into the bound we get

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

1015 where $c_{0,\min} := \min_{j \in [d]} c_{0,j}$. If we choose $K = \mathcal{O}(\varepsilon^{-4})$ we get that

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

1016 C.2 Convergence under PL-condition

1017 **Theorem C.2.** Let f satisfies PL-condition and each f_i satisfies Definition B.1. Assume that
1018 Assumption B.4 holds. Then the iterates of NGN-D (Alg. 3) with step-size parameters $\{c_j\}_{j=1}^d$ such
1019 that $c_j \leq \min\{1/2L_j, 6/\mu\}$ satisfy

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

1020 where $c_{\min} := \min_{j \in [d]} c_j$. Moreover, if $c_j = \mathcal{O}(\varepsilon)$ for all $j \in [d]$ then after $K =$
1021 $\max\{\mathcal{O}(\varepsilon^{-1}), \mathcal{O}(1)\} \log \varepsilon^{-1}$ iterations we obtain $\mathbb{E} [f(x^K) - f^*] \leq \mathcal{O}(\varepsilon)$.

1022 *Proof.* We obtain (19) and use Definition B.3

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

1023 Subtracting f^* from both sides of the inequality above and taking full expectation we obtain

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

1024 Unrolling the recursion above for $\{0, \dots, K-1\}$ iterations we derive

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

1025 Now we follow the proof of Lemma A.3 in Garrigos and Gower [20]. Let us choose $c_j =$
 1026 $\min\{1/2L_j, \varepsilon/2dA_j\}$. Together with the choice of $K \geq \max_{j \in [d]} \max \left\{ \frac{1}{\varepsilon} \frac{12A_j}{\mu}, \frac{12L_j}{\mu} \right\} \log \frac{2(f(x^0) - f^*)}{\varepsilon}$
 1027 we get

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

1028 Now we have two cases:

1029 1. c_{\min} does not depend on ε , then we have

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

1030 2. c_{\min} does depend on ε , i.e. $c_{\min} = \mathcal{O}(\varepsilon)$, then we have

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

1031 Therefore, combining all together we get

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

1032 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.

1033 □

1034 D Convergence of NGN-M

1035 **Theorem 4.3.** Let Assumptions 4.1, 4.2 hold. Let the step-size hyperparameter $c > 0$ and the
 1036 momentum parameter $\beta = \frac{\lambda}{1+\lambda}$ be constants where $\lambda \leq \min\{cL, 0.5(1+cL)^{-1}(1+2cL)^{-1}\}$. Then
 1037 the iterates of NGN-M (Alg. 1) satisfy

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

1038 where \bar{x}^{K-1} is chosen uniformly at random from $\{x^0, \dots, x^{K-1}\}$. Moreover, if we set $c = \mathcal{O}(1/\sqrt{K})$
 1039 then we obtain $\mathbb{E} [f(\bar{x}^{K-1}) - f(x^*)] \leq \mathcal{O}(1/\sqrt{K})$.

1040 **Remark D.1.** In fact, if $\lambda \leq \frac{1}{(1+cL)(1+2cL)}$, then it implies that $\lambda \leq \frac{1}{cL}$ because $\frac{1}{x} > \frac{1}{(1+x)(1+2x)}$
 1041 for any $x > 0$.

1042 *Proof.* To prove the convergence of NGN-M we consider IMA formulation Equation (5):

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

1043 where $\beta = \frac{\lambda}{1+\lambda}$, $z^{k+1} = x^{k+1} + \lambda(x^{k+1} - x^k)$.

1044 At iteration $k = 0$ we have

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

1045 Therefore, we get

$$\begin{aligned} \|z^1 - x^*\|^2 &= \|z^0 - x^*\|^2 - 2\gamma_0 \langle \nabla f_{S_0}(x^0), z^0 - x^* \rangle + \gamma_0^2 \|\nabla f_{S_0}(x^0)\|^2 \\ &\stackrel{\text{Lem. B.6}}{\leq} \|z^0 - x^*\|^2 - 2\gamma_0 \langle \nabla f_{S_0}(x^0), x^0 - x^* \rangle + \frac{4cL}{1+2cL} \gamma_0 (f_{S_0}(x^0) - f_{S_0}^*) \\ &\quad + \frac{2c^2L}{1+cL} \max \left\{ \frac{2cL-1}{2cL+1}, 0 \right\} f_{S_0}^*. \end{aligned} \tag{20}$$

1046 Let $\gamma_0 = \rho + \tilde{\gamma}_0$ where $\rho = \frac{c}{(1+cL)(1+2cL)}$. Then we have

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

1047 Using the above we continue from (20)

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

1048 Here we have

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

1049 $\tilde{\gamma}_0 \leq \frac{3c^2L}{1+2cL}$, and $f_{S_0}(x^0) - f_{S_0}^* \geq 0$. Hence, we get

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

1050 Rearranging terms and taking expectation we get

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

1051 Next, for $k > 0$ we can use the relation $z^k = x^k + \lambda(x^k - x^{k-1})$. We expand $\|z^{k+1} - x^*\|^2$

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

1052 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
1053 defined later. Therefore, we have

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

1054 We need to find ρ such that

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

1055 Since $\tilde{\gamma}_k = \gamma_k - \rho$, then we have

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

1056 The inequality above is satisfied if it is satisfied for the lower bound on γ_k (which is $c/(1+cL)$), i.e.

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

1057 We can take $\rho = \frac{c}{(1+cL)(1+2cL)}$ since $\lambda \geq 0$.

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

1058 Using the above, we get from (23)

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

1059 Taking expectations we get

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

1060 Rearranging terms we get

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

1061 Combining Equation (22) and Equation (25) for iterations $\{1, \dots, K-1\}$ we get

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

1062 We need to ensure that $\rho - c\lambda > 0$ which is satisfied for λ such that

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

1063 Note that we also assume that $\lambda \leq cL$. Therefore, from (26) we get

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

1064 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

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

1065 where we use the convexity of f and Jensen's inequality. Plugging the value of $\rho = \frac{c}{(1+cL)(1+2cL)}$
1066 inside we get

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

1067 Choosing $c = \mathcal{O}(1/\sqrt{K})$ we get

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

1068 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
1069 derive the statement of the theorem. \square

1070 E Convergence of NGN-M with Decaying Step-size

1071 **Lemma E.1.** *We have*

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

1072 *Proof.* We refer to Lemma A.8 from Garrigos and Gower [20]. \square

1073 To prove the convergence of NGN-M with decaying c_k we consider IMA formulation (see Section A
1074 in the paper):

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

1075 where $c_k = \frac{c_0}{\sqrt{k+1}}$, $\lambda_k = Lc_k$, $\lambda_0 = 0$.

1076 **Theorem E.2.** *Assume that each f_i is convex and L -smooth, and that Assumption 3.2 holds. Let the
1077 step-size hyperparameter is set $c_k = \frac{c_0}{\sqrt{k}}$, momentum parameter $\lambda_k \leq \min\{c_k L, 0.5(1 + c_k L)^{-1}(1 +$
1078 $2c_k L)^{-1}\}$. Then the iterates of NGN-M satisfy*

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

1079 where $\hat{x}^{K-1} = \sum_{k=0}^{K-1} \frac{\rho_k}{\sum_{k=0}^{K-1} \rho_k} x^k$, $\rho_k = \frac{c_k}{(1+c_k L)(1+2c_k L)}$.

1080 *Proof.* At iteration $k = 0$ we have

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

1081 Therefore, we get

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

1082 Let $\gamma_0 = \rho_0 + \tilde{\gamma}_0$ where $\rho_0 = \frac{c_0}{(1+c_0L)(1+2c_0L)}$. Then we have

$$\begin{aligned}
\tilde{\gamma}_0 &= \gamma_0 - \rho_0 \\
&\stackrel{\text{Lem. B.5}}{\leq} c_0 - \frac{c_0}{(1+c_0L)(1+2c_0L)} \\
&= c_0 \frac{1+3c_0L+2c_0^2L^2-1}{(1+c_0L)(1+2c_0L)} \\
&= c_0^2L \frac{3+3c_0L}{(1+c_0L)(1+2c_0L)} \\
&= \frac{3c_0^2L}{1+2c_0L}.
\end{aligned}$$

1083 Using the above we continue from (33)

$$\begin{aligned}
\|z^1 - x^*\|^2 &\stackrel{\text{conv.}}{\leq} \|z^0 - x^*\|^2 - 2\gamma_0(f_{S_0}(x^0) - f_{S_0}(x^*)) + \frac{4c_0L}{1+2c_0L}\gamma_0(f_{S_0}(x^0) - f_{S_0}^*) \\
&\quad + \frac{2c_0^2L}{1+c_0L} \max\left\{\frac{2c_0L-1}{2c_0L+1}, 0\right\} f_{S_0}^* \\
&\leq \|z^0 - x^*\|^2 - 2\rho_0(f_{S_0}(x^0) - f_{S_0}(x^*)) - 2\tilde{\gamma}_0(f_{S_0}(x^0) - f_{S_0}^*) + 2\tilde{\gamma}_0(f_{S_0}(x^*) - f_{S_0}^*) \\
&\quad + \frac{4c_0L}{1+2c_0L}\gamma_0(f_{S_0}(x^0) - f_{S_0}^*) + \frac{2c_0^2L}{1+c_0L} \max\left\{\frac{2c_0L-1}{2c_0L+1}, 0\right\} f_{S_0}^* \\
&= \|z^0 - x^*\|^2 - 2\rho_0(f_{S_0}(x^0) - f_{S_0}(x^*)) - 2\left(\gamma_0 - \rho_0 - \frac{2c_0L}{1+2c_0L}\gamma_0\right)(f_{S_0}(x^0) - f_{S_0}^*) \\
&\quad + 2\tilde{\gamma}_0(f_{S_0}(x^*) - f_{S_0}^*) + \frac{2c_0^2L}{1+c_0L} \max\left\{\frac{2c_0L-1}{2c_0L+1}, 0\right\} f_{S_0}^*. \tag{34}
\end{aligned}$$

1084 Here we have

$$\begin{aligned}
\gamma_0 - \rho_0 - \frac{2c_0L}{1+2c_0L}\gamma_0 &= \frac{1}{1+2c_0L}\gamma_0 - \rho_0 \\
&= \frac{1}{1+2c_0L}\gamma_0 - \frac{c_0}{(1+c_0L)(1+2c_0L)} \\
&\stackrel{\text{Lem. B.5}}{\geq} \frac{1}{1+2c_0L} \frac{c_0}{1+c_0L} - \frac{c_0}{(1+c_0L)(1+2c_0L)} \\
&= 0,
\end{aligned}$$

1085 $\tilde{\gamma}_0 \leq \frac{3c_0^2L}{1+2c_0L}$, and $f_{S_0}(x^0) - f_{S_0}^* \geq 0$. Hence, we get

$$\begin{aligned}
\|z^1 - x^*\|^2 &\leq \|z^0 - x^*\|^2 - 2\rho_0(f_{S_0}(x^0) - f_{S_0}(x^*)) + \frac{6c_0^2L}{1+2c_0L}(f_{S_0}(x^*) - f_{S_0}^*) \\
&\quad + \frac{2c_0^2L}{1+c_0L} \max\left\{\frac{2c_0L-1}{2c_0L+1}, 0\right\} f_{S_0}^*.
\end{aligned}$$

1086 Rearranging terms and taking the expectation we get

$$\begin{aligned}
2\rho_0\mathbb{E}[f(x^0) - f(x^*)] &\leq \mathbb{E}[\|z^1 - x^*\|^2] - \|z^0 - x^*\|^2 + \frac{6c_0^2L}{1+2c_0L}\sigma_{\text{int}}^2 \\
&\quad + \frac{2c_0^2L}{1+c_0L} \max\left\{\frac{2c_0L-1}{2c_0L+1}, 0\right\} \sigma_{\text{pos}}^2. \tag{35}
\end{aligned}$$

1087 Next, for $k > 0$ we can use the relation $z^k = x^k + \lambda_k(x^k - x^{k-1})$. We expand $\|z^{k+1} - x^*\|^2$

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

1088 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
1089 defined later. Therefore, we have

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

1090 We need to find ρ_k such that

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

1091 Since $\tilde{\gamma}_k = \gamma_k - \rho_k$, then we have

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

1092 The inequality above is satisfied if it is satisfied for the lower bound on γ_k (which is $c/(1+cL)$), i.e.

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

1093 We can take $\rho_k = \frac{c_k}{(1+c_k L)(1+2c_k L)}$ since $\lambda \geq 0$.

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

1094 Using the above, we get from (36)

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

1095 Taking expectations, we get

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

1096 Rearranging terms, we get

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

1097 Combining (35) and (38) for iterations $\{1, \dots, K-1\}$ we get

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

1098 Note that choosing $\lambda_k = \min\{c_kL, 0.5(1+c_kL)^{-1}(1+2c_kL)^{-1}\}$ ensures that $\frac{\rho_k}{2} \geq c_k\lambda_k$. In-
1099 deed, we have

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

1100 Therefore, from (39) and the facts that $\lambda_0 = 0$ and $\lambda_k \leq c_kL$ we get

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

1101 We have by Lemma E.1

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

1102 Therefore, using (41), $z^0 = x^0$ in (40) and dividing both sides in (40) by $\sum_{k=0}^{K-1} \rho_k$ we derive

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

1103 With an lower bound on $\sum_{k=0}^{K-1}$ and Jensen's inequality we conclude that

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

1104 where $\hat{x}^{K-1} = \sum_{k=0}^{K-1} \frac{\rho_k}{\sum_{k=0}^{K-1} \rho_k} x^k$.

1105

□

1106 F Stability of NGN-M on a Simple Problem

1107 We consider 1D convex functions of the form $f(x) = Lx^2(1+p^2(x))$ that satisfy the following
1108 assumption.

1109 **Assumption F.1.** There exists a constant C such that $C(1+p^2(x)) \geq xp(x)p'(x)$.

1110 Note that $1+p^2(x) \geq 1$ and $\deg(1+p^2(x)) = \deg(xp(x)p'(x))$. Therefore, this assumption is mild.

1111 *Remark F.2.* For example, the function $f(x) = x^2(1+x^2)$ (i.e., $p(x) = x$) is convex and satisfies
1112 Assumption F.1 with $C = 1$.

1113 *Remark F.3.* Let $p(x) = \sum_{j=0}^m a_j x^j$. Then for large values of x in magnitude, $p(x) \sim$
1114 $a_m x^m$, $p'(x) \sim m a_m x^{m-1}$. Therefore, the constant C should be expected of order $C \approx m$, where
1115 $m = \deg(p(x))$.

1116 The function $f(x)$ is non-negative for any $x \in \mathbb{R}$ and its minimum $f^* = 0$ is attained at $x = 0$ by
1117 design. Let us compute a step of NGN-M on this problem

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

1118 Note that the convexity of f implies that

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

1119 In particular, (45) implies that $1 + p^2(x) + xp(x)p'(x) \geq \frac{1}{2}(1 + p^2(x)) > 0$. Therefore, we can
1120 obtain lower and upper bounds on $\hat{\gamma}_k$.

1121 **Lemma F.4.** *Let Assumption F.1 hold with a constant $C > 0$ and $f(x) = x^2(1 + p^2(x))$ be convex.*

1122 *Let $c \geq \frac{1}{2L}$. Then we have $\hat{\gamma}_k \in \left[\frac{1}{2(1+C)}, 2\right]$.*

1123 *Proof.* Indeed, the upper bound on $\hat{\gamma}_k$ follows from the following inequality

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

1124 due to (45). The lower bound can be obtained as follows

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

1125 \square

1126 The update rule of NGN-M can be rewritten as

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

1127 Let us consider the joint dynamics of $w^k := ([x^k]^\top, [x^{k-1}]^\top)^\top \in \mathbb{R}^{2d}$. We have that

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

1128 where

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

1129 Now we are ready to prove the convergence of NGN-M on this simple problem for any value $c \geq \frac{1}{2}$.

1130 **Theorem F.5.** Let $f(x) = x^2(1 + p^2(x))$ be convex and Assumption F.1 holds. Let $\beta \geq \frac{(2(1+C)-1)^2}{(2(1+C)+1)^2}$
 1131 and $c \geq \frac{1}{2L}$. Then the iterates of NGN-M on $f(x)$ converge to the minimum $f^* = 0$.

1132 *Proof.* We follow the standard proof of SGD with Polyak momentum [59]. At this stage, we need
 1133 to estimate the eigenvalues of \mathbf{G} . To do so, we will proceed with a permutation matrix Π^4 which
 1134 transforms the matrix \mathbf{G} to the block-diagonal matrix as

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

1135 where

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

1136 Since the matrix \mathbf{G} is a block-diagonal matrix, we have $\|\mathbf{G}\| \leq \max_i \|\mathbf{G}_i\|$. Therefore, the problem
 1137 is now simplified to bounding the spectral radii of the individual blocks \mathbf{G}_i , for $i = 1, 2, \dots, d$. The
 1138 two eigenvalues u_1 and u_2 of \mathbf{G}_i are the roots of the quadratic

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

1139 which take different values depending on the discriminant $\Delta := (1 + \beta - (1 - \beta)\hat{\gamma}_k)^2 - 4\beta$. Let us
 1140 find the values of β when the discriminant is negative. We need to satisfy the inequality

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

1141 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
 1142 we satisfy the last inequality, and consequently the discriminant is non-positive, if we choose

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

1143 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
 1144 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
 1145 equal in absolute value

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

1146 This gives us that $\|\mathbf{G}_i\| \leq \sqrt{\beta} < 1$. Therefore, the algorithm converges for any value of β in this
 1147 range.

1148 It remains to use Lemma 11 from Foucart [18] which says that for a given matrix $\mathbf{A} \in \mathbb{R}^{d \times d}$, and
 1149 $\epsilon > 0$, there exists a matrix norm $\|\cdot\|$ such that

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

1150 where $\rho(\mathbf{A}) = \max\{|\lambda| : \lambda \text{ eigenvalue of } \mathbf{A}\}$ (spectral radius of \mathbf{A}).

1151 Asymptotically ⁵ (as $k \rightarrow \infty$, one can show (see Theorem 12 in [18]) that

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

1152 where $\rho(\mathbf{G}) \leq \sqrt{\beta} < 1$ in our analysis. Therefore, NGN-M with hyperparameters $c \geq \frac{1}{2}$ and $\beta \geq \frac{1}{0}$
 1153 converges. □

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.

⁵A non-asymptotic version of the analysis can be derived using Theorem 5 by [77]

1154 *Remark F.6.* For example, NGN-M converges on $f(x) = x^2(1 + x^2)$ for any $c \geq \frac{1}{2}$ and $\beta \geq \frac{9}{25}$.

1155 Theorem F.5 shows that NGN-M remains stable even with an arbitrarily large step-size hyperparameter
 1156 c . Thanks to the adaptive nature of NGN step-size, the actual update scale is automatically shrunk
 1157 when necessary, preserving convergence. Importantly, this is possible with a choice of momentum
 1158 parameter β close to 1, which extends the results of Section 4. We acknowledge that our current
 1159 analysis is restricted to the special convex class of 1D functions $f(x) = x^2(1 + p^2(x))$ satisfying
 1160 Assumption F.1. Extending such stability guarantees to wider function classes with large momentum
 1161 β remains a significant open challenge.

1162 To support the theoretical result, we test the performance of NGN-M and GDM (Gradient Descent with
 1163 Momentum) on the problem $f(x) = x^2(1 + x^2)$, which is convex and satisfies Assumption F.1; see
 1164 Figure F.1. We run both algorithms, varying the step-size hyperparameter in $\{10^{-4}, \dots, 10^4\}$. We run
 1165 algorithms for 10^5 iterations. We stop training if the loss reaches a threshold 10^{-15} or exceeds 10^{10}
 1166 for the first time. We observe that (i) for small step-size hyperparameters, both methods converge
 1167 but do not reach the threshold 10^{-15} ; (ii) NGN-M reaches the threshold even for extremely large
 1168 values of the step-size hyperparameter while GDM diverges. (iii) the fastest convergence of GDM is
 1169 achieved with the step-size hyperparameter 10^{-2} after 691 iterations while the fastest convergence
 1170 of NGN-M is achieved with $c = 10^1$ after 269 iterations. In other details, NGN-M achieves faster
 1171 convergence and much more stable to the choice of the step-size hyperparameter. These results align
 1172 well with our theoretical analysis.

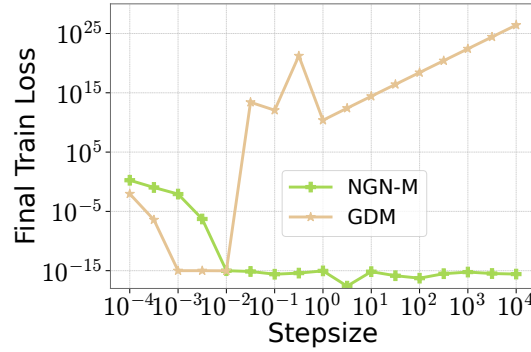


Figure F.1: Comparison of SGDM and NGN-M when minimizing a function $f(x) = x^2 + x^4$.

1173 G How to Derive Diagonal NGN-based Step-size?

1174 Here we provide derivations of how combine NGN and diagonal step-size following Section 3.3 for
1175 completeness.

1176 We consider the following model

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

1177 where $r(x) = \sqrt{f(x)}$. We compute the gradient of RHS of (59) w.r.t. p and equal it to zero:

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

1178 Therefore, we have

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

1179 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

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

1180 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

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

1181 G.1 Design Comparison of NGN-MDv1 and NGN-MDv2

1182 The derivations in (3) are used to provide an intuition of how one can add a diagonal step-size into
1183 NGN by choosing the regularization matrix Σ_k . By choosing $\Sigma_k = \mathbf{D}_k$ we recover the update
1184 direction of NGN-MDv1. In this case, we have only one global NGN step-size in front of \mathbf{D}_k . The
1185 design of NGN-MDv2 follows a more straightforward intuition. In particular, it can be seen as a direct
1186 extension of NGN to diagonal case by replacing the squared gradient norm $\|\nabla f_{S_k}(x^k)\|^2$ by the
1187 squared partial derivative $(\nabla_j f_{S_k}(x^k))^2$ for each parameter $j \in [d]$.

1188 The main difference in comparison with Adam is the order in which the preconditioning and mo-
1189 mentum is applied. In both NGN-MDv1 and NGN-MDv2 we average the preconditioned updates
1190 $\Sigma_k^{-1} \nabla f_{S_k}(x^k)$, i.e. we first apply preconditioning and momentum later. In contrast, in Adam the
1191 stochastic gradients are averaged to construct new momentum term, and then the momentum is pre-
1192 conditioned. In other words, the momentum is applied first and then it is followed by preconditioning.
1193 We believe this change might be one of the reasons behind the step-size hyperparameter resilience as
1194 well.

1195 In practice, we found out that the tuned performance of NGN-MDv1 is slightly better than that of
1196 NGN-MDv2. Moreover, NGN-MDv1 demonstrates higher resilience to the choice of the step-size
1197 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 [42], 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 [69] or ZeRO-2 [62], 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 [84]. Despite its widespread application, the role of weight decay is poorly understood. Andriushchenko et al. [1] 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 [88].

1238 We propose two ways of adding weight decay to NGN-MDv1. The first variant follows the approach
 1239 of [46], adding decoupled weight decay λ :

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

1240 In this update rule, the weight is added separately from the update direction $\Sigma_k^{-1} \nabla f_{S_k}(x^k)$. We call
 1241 the resulting algorithm (60) Dec-NGN-MDv1, that stands for decoupled NGN-MDv1.

1242 H.1 Combining NGN-MDv1 and Weight Decay Regularization

1243 We now discuss how to combine NGN-MDv1 and weight decay, following the idea that weight decay
 1244 should perform weight regularization.

1245 We consider the following model

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

1246 By taking the gradient of $f_{\Sigma_k, \lambda}$ w.r.t. p we get

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

1247 Therefore, we get

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

1248 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
 1249 $u = v = \sqrt{2} \nabla r(x^k)$ we get that

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

1250 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 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 $\max\{0, \cdot\}$.
-

1251 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

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

1252 To summarize, the update of NGN-Dv1W is the following

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

1253 To prevent the step-size next to $\Sigma_k^{-1} \nabla f(x^k)$ from being negative, the final update has the form

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

1254 where $[\cdot]_+ := \max\{\cdot, 0\}$. Now we can add momentum on top and obtain the following update of
1255 NGN-MDv1W

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

1256 This combination of NGN-MDv1 and weight decay is summarized in Algorithm 4. We highlight
1257 that now the weight decay is incorporated inside the adaptive step-size as well as regularizing the
1258 coefficient next to x^k .

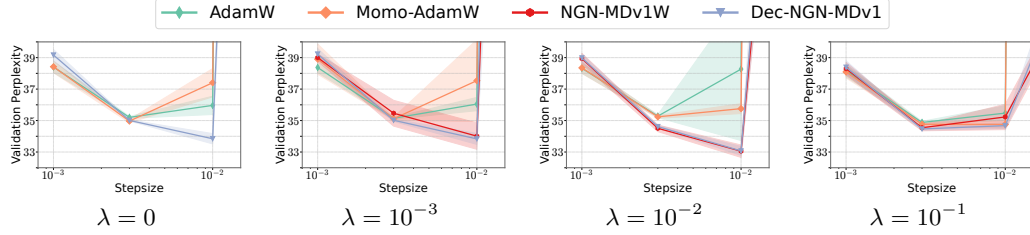


Figure H.1: Adding weight decay when pretraining a 70M Transformer++. When properly tuned, a value of weight decay > 0 enhances the performance of all algorithms. NGN-MDv1 retains his characteristic stability, and achieves smaller perplexity in all scenarios.

1259 H.2 Empirical Validation of the Proposed Combinations

1260 Having two possible ways of adding weight decay to NGN-MDv1, we test them on pretraining a 70M
 1261 transformer on language modeling. The validation perplexity at the end of training is reported in
 1262 Figure H.1. We note that when weight decay is turned off, both NGN-MDv1W and Dec-NGN-MDv1
 1263 reduce to NGN-MDv1.

1264 First, we observe that when weight decay is properly tuned, all algorithms improve over the baseline
 1265 case with no weight decay, which is consistent with the observation of Xiao [84] and Andriushchenko
 1266 et al. [1] on AdamW. We also note that Dec-NGN-MDv1 and NGN-MDv1W require a smaller weight
 1267 decay value compared to the other algorithms. Finally, the stability and performance of NGNMDv1 are
 1268 preserved by both variations, allowing training with larger learning rates, and significantly improving
 1269 over AdamW and Momo-Adam.

1270 We do not observe a substantial difference between the two proposed modifications of NGN-MDv1
 1271 for this task. We remark however that these two versions serve substantially different purposes, and
 1272 pretraining language models might not be the most representative task to evaluate the effect of adding
 1273 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.

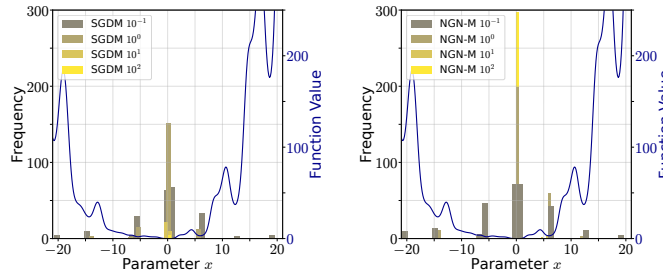


Figure I.1: Comparison of SGDM and NGN-M when minimizing function in (64).

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.

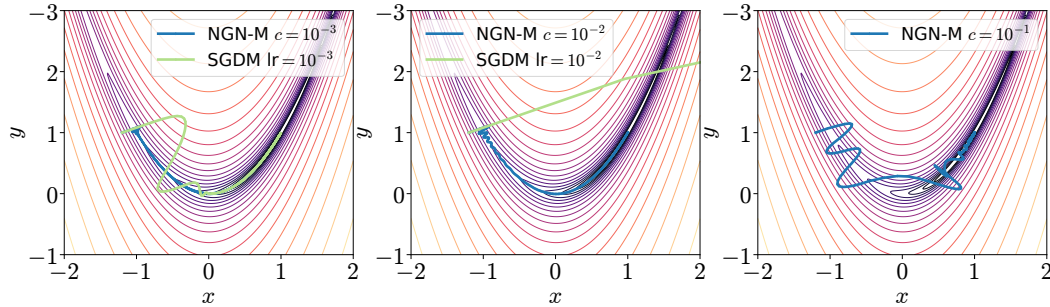


Figure I.2: Trajectories of NGN-M and SGDM when minimizing the Rosenbrock function and varying the step-size hyperparameter.

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

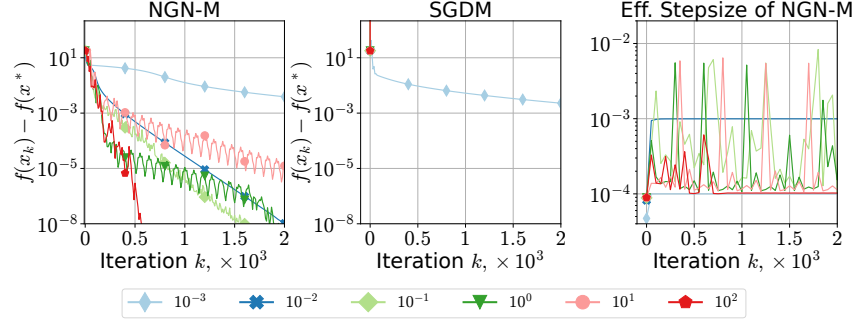


Figure I.3: Training dynamics of NGN-M and SGDM when minimizing the Rosenbrock function and varying the step-size hyperparameter.

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 4.3 and Theorem E.2 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 .

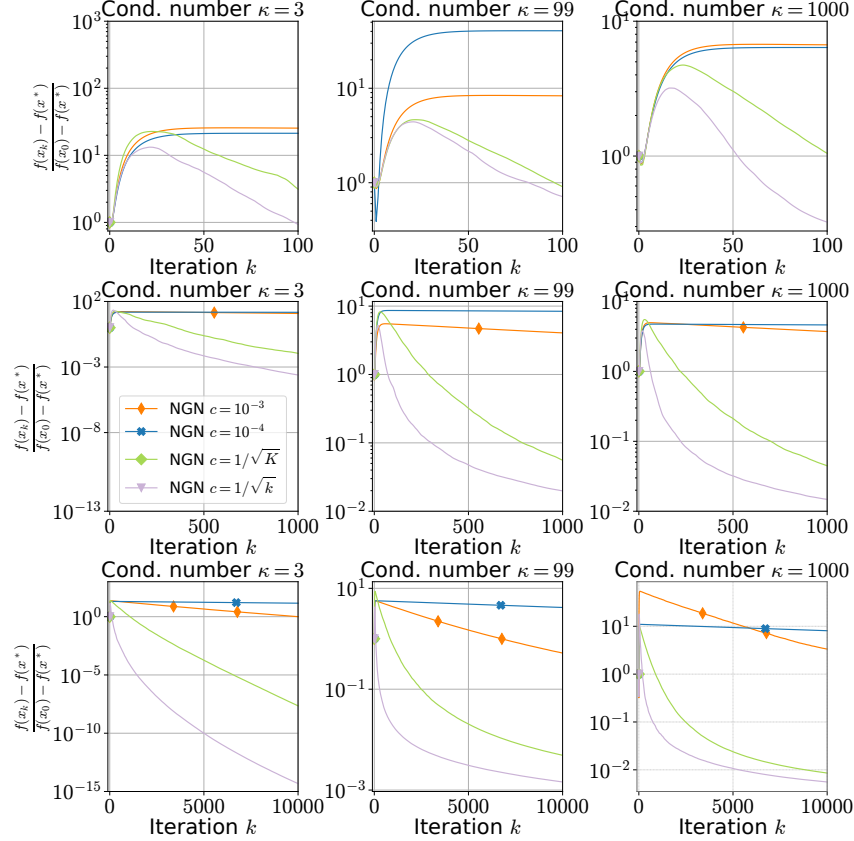


Figure I.4: Training dynamics of NGN-M with several choices of the step-size hyperparameter varying the condition number of the quadratic problem.

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 [58] 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. [66], and use Pytorch Image Models codebase [81]. 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. [66], we train in `bf16`, instead of `float16`, and do not employ weight decay regularization.

For pre-training Transformers on Causal Language Modeling, we build upon the nanoGPT [35] implementation, augmenting it with Rotational Positional Embedding [73], RMSNorm [87], and SwiGLU [68]. We call this enhanced version Transformer++. Models are trained with a batch size of 256, context length of 2048 tokens, vocabulary size of 50280 and make use of GPT-Neox tokenizer [4]. We adopt an enhanced training recipe, made popular by large language models such as LLaMa [75]. These modifications include: training in `bf16`; 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

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	Wikitext-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
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

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

1331 tying between embedding and last linear layer. All models are trained on SlimPajama-627B [72], a
1332 cleaned and deduplicated version of RedPajama. We report validation perplexity on a separate subset
1333 of Slim-Pajama consisting of 10M tokens. The total compute is estimated following Kaplan et al.
1334 [33], where the estimated number of floating-point operations (FLOPs) is $6 \times \text{Number of Parameters}$
1335 $\times \text{Number of Tokens}$.

1336 Experiments of small and middle size are performed on 1xRTX 4090. We perform ImageNet32
1337 experiments on 2xA100-40GB, and ImageNet1k experiments on 4xA100-SXM4-40GB. For pretrain-
1338 ing Transformers on Language Modeling, we employ 8xH100-HBM3-80GB GPUs. With multiple
1339 devices in use, we employ Distributed Data Parallel to parallelize the training process.

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

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 [70], 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 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 [26], PTB [50]), and (Transformer [35], Tiny Shakespeare [34]) 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 [49]) and (Transformer, Rotten Tomatoes [57]) 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

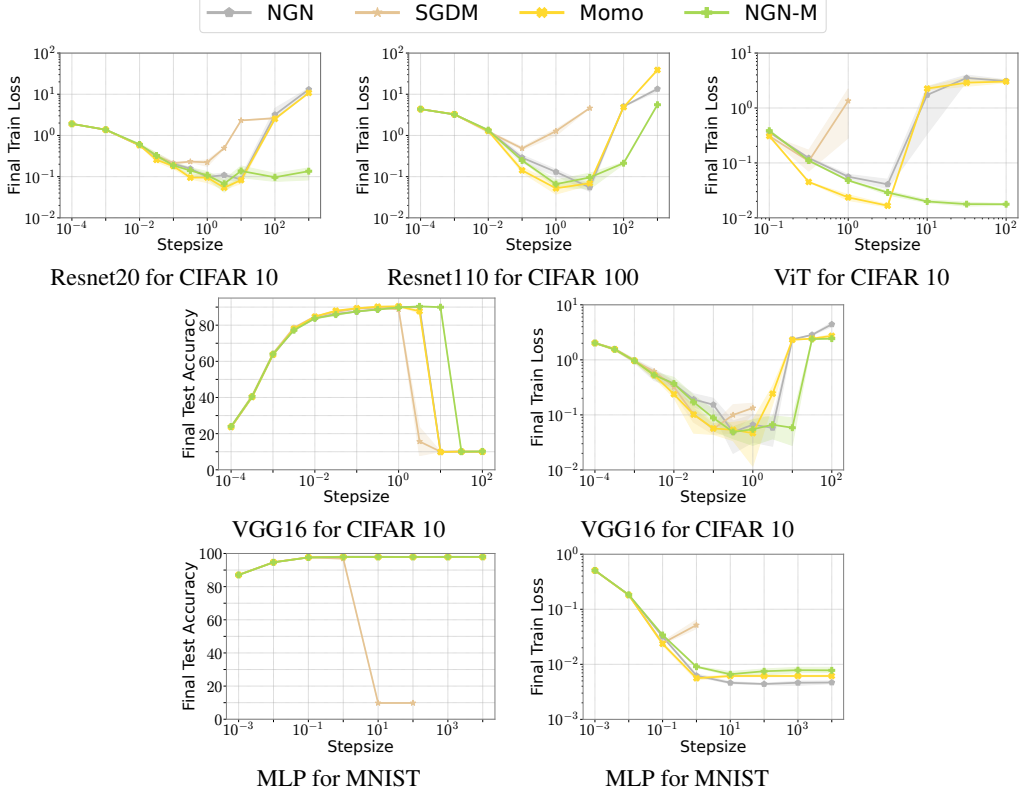


Figure J.1: Stability performance of algorithms supporting momentum varying step-size hyperparameter (c for NGN and NGN-M, α_0 for Momo, and step-size for SGDM). We observe that NGN-M achieves the training loss close to the best possible for a wider range of the step-size hyperparameter.

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 [8], Adabound [47], and Adabelief [94]. The results are presented in Table 5.

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.

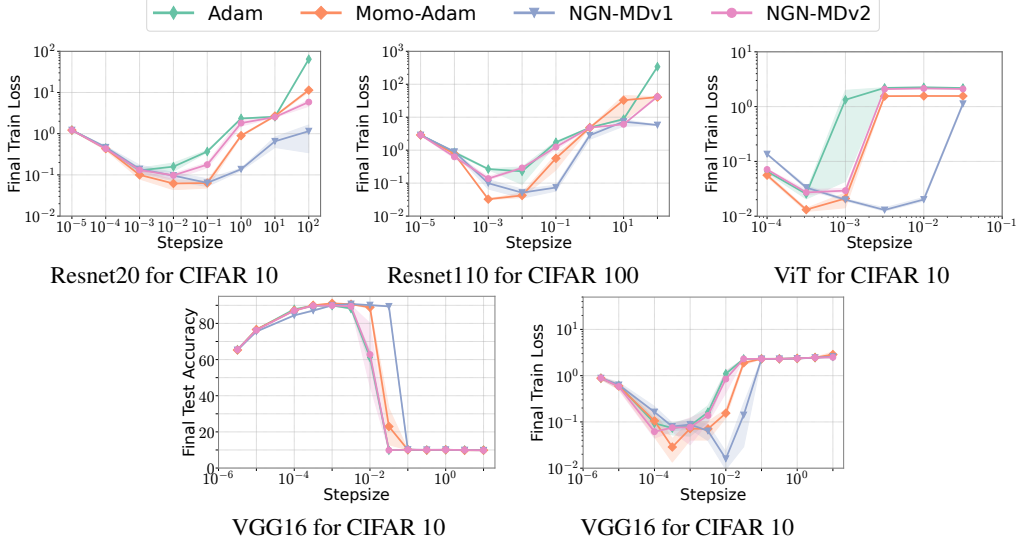


Figure J.2: Stability performance of algorithms supporting momentum and diagonal step-size varying step-size hyperparameter (c for NGN-MDv1 and NGN-MDv2, α_0 for Momo-Adam, and step-size for Adam). We observe that NGN-MDv1 achieves the training loss close to the best possible for a wider range of the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

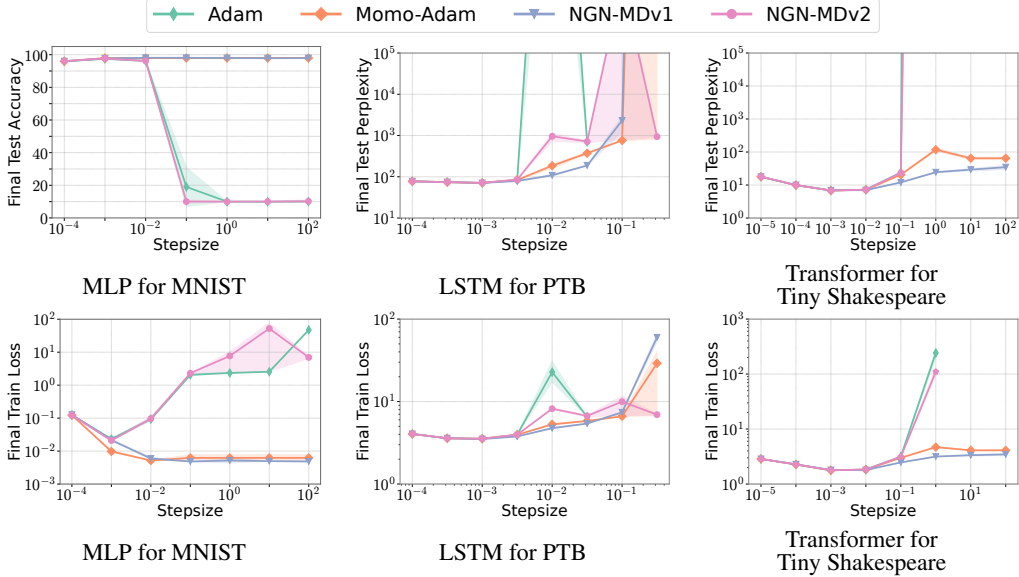


Figure J.3: Stability performance of algorithms supporting momentum and diagonal step-size varying step-size hyperparameter (c for NGN-MDv1 and NGN-MDv2, α_0 for Momo-Adam, and step-size for Adam). We observe that NGN-MDv1 achieves the training loss close to the best possible for a wider range of the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

1388 Since the adaptivity of Adam is only in the second-order momentum applied as a normalization, in
1389 our experiment we compare the following quantities

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

1390 where γ is the step-size hyperparameter of Adam.

1391 Let us first describe the results for ResNet20 in Figure J.6. We observe that NGN-MDv1 tends to set
1392 smaller effective step-size compared to two other algorithms. This is especially visible for the large

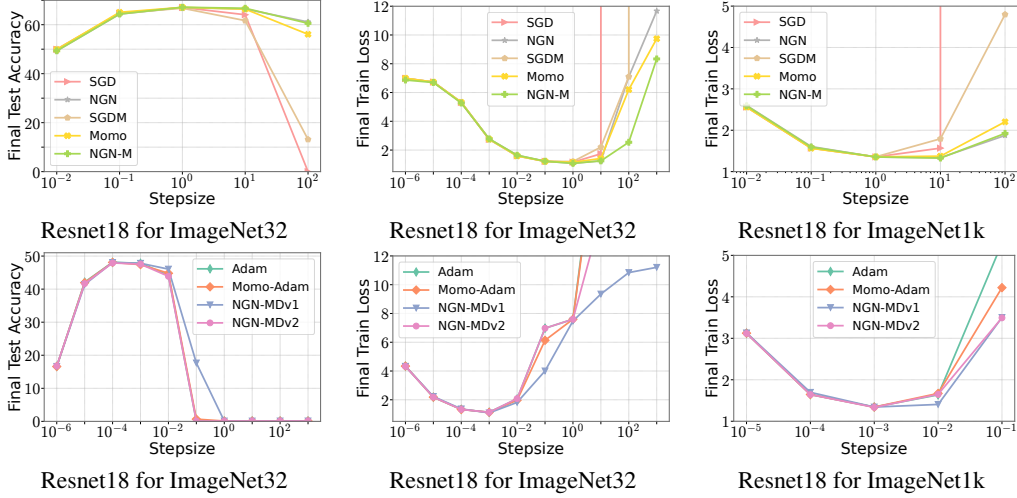


Figure J.4: Stability performance of algorithms supporting momentum (**first row**), and momentum with diagonal step-size (**second row**) varying step-size hyperparameter (c for NGN, NGN-M, NGN-MDv1, and NGN-MDv2, α_0 for Momo and Momo-Adam, and step-size for SGD, SGDM, and Adam).

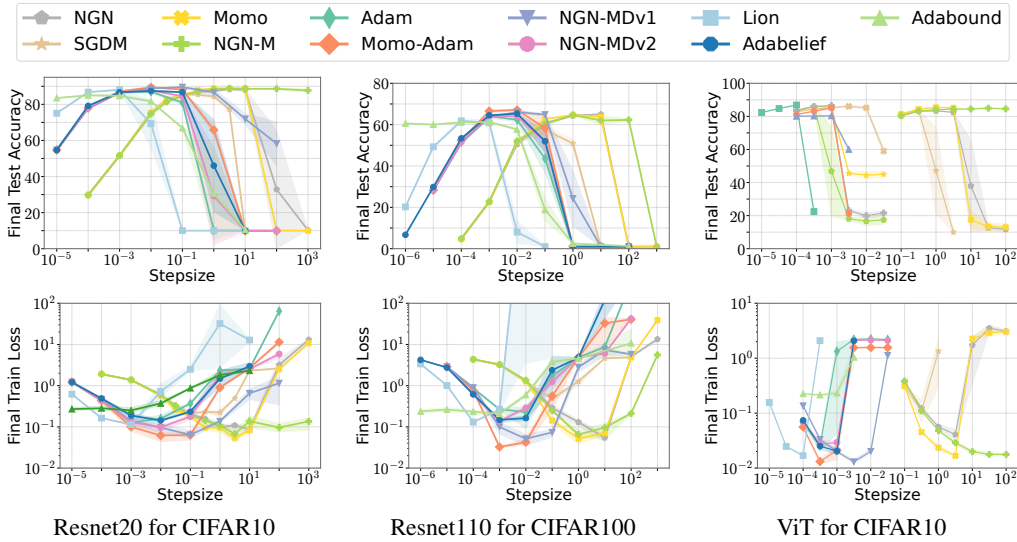


Figure J.5: Stability performance of various optimizers for (Resnet20, CIFAR10), (Resnet110, CIFAR100), (ViT, CIFAR10) workloads.

1393 step-size hyperparameter values where the adaptive step-size of NGN-MDv1 is by several orders in
 1394 magnitude smaller than that of Adam and Momo-Adam. In contrast, the coordinate-wise adaptive
 1395 step-size of Momo-Adam is mostly follow that of Adam. Considering that the stability performance
 1396 of NGN-MDv1 is much higher for this task, this happens mainly due to the fact that the adaptation
 1397 mechanism of NGN-MDv1 step-size is more conservative than that of Momo-Adam.

1398 Now we switch to the results on ViT model in Figure J.7. Here both Momo-Adam and NGN-MDv1
 1399 tend to utilize smaller effective coordinate-wise step-size, by several orders in magnitude smaller
 1400 than that of Adam. However, the adaptation mechanism of NGN-MDv1 is still more conservative
 1401 than that of Momo-Adam, especially for large step-size hyperparameters. We also highlight that in
 1402 this experiment the best performance of NGN-MDv1 is achieved with $c = 10^{-3}$. When we vary the
 1403 step-size hyperparameter c , the effective coordinate-wise step-size does not change dramatically,
 1404 especially for layers.0.0.fin.to_qkv layer.

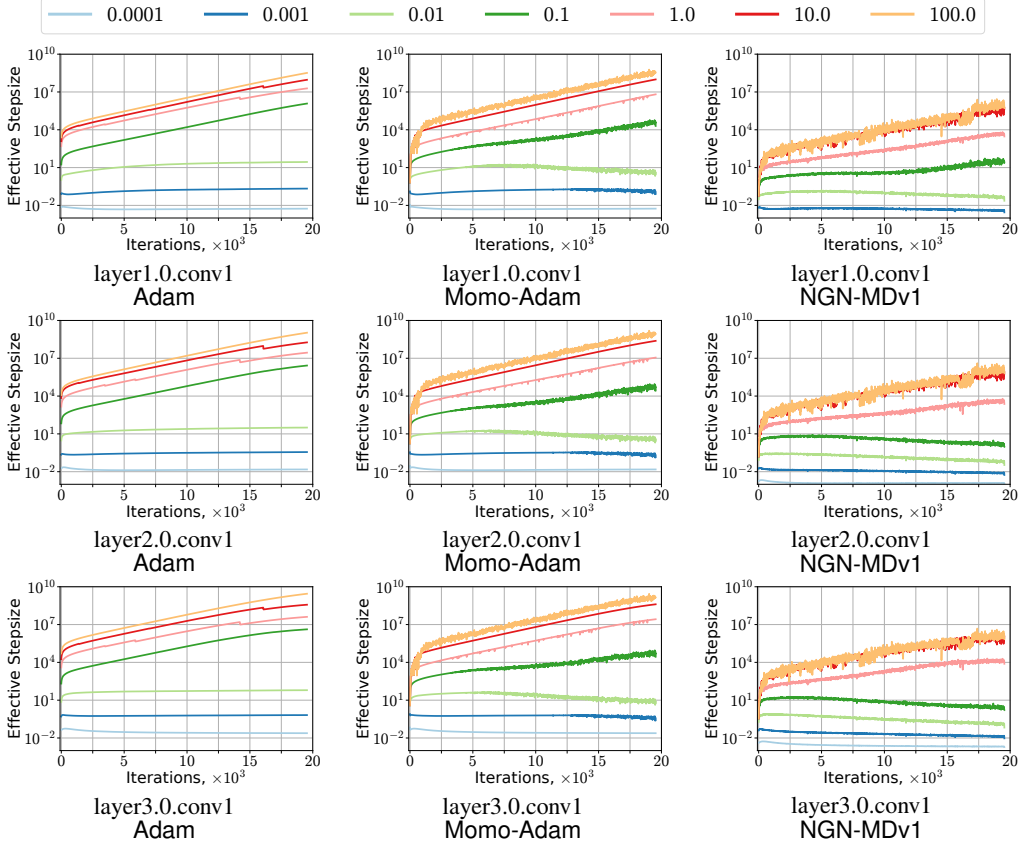


Figure J.6: The adaptive stepsize of Adam (**first column**), Momo-Adam (**second column**), and NGN-MDv1 (**third column**) algorithms in training ResNet20 model on CIFAR10 dataset. We plot the average stepsize $\frac{\gamma}{(\mathbf{D}_k)_{(j)}}$ (for Adam), $\frac{\tau_k}{(\mathbf{D}_k)_{(j)}}$ (for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

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.

J.8 Comparison of Algorithms with Diagonal Step-size

Now we compare algorithms with diagonal step-size such as NGN-D, Adagrad [16], and RMSprop [36]. 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

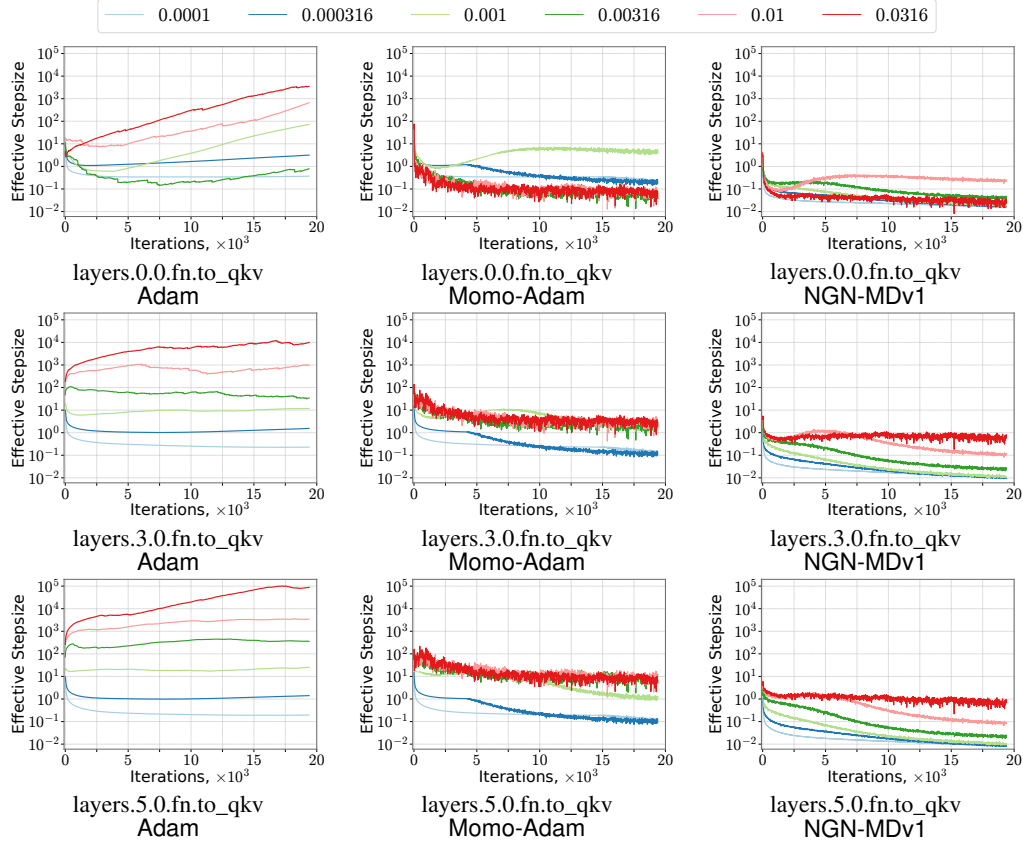


Figure J.7: The adaptive stepsize of Adam (**first column**), Momo-Adam (**second column**), and NGN-MDv1 (**third column**) algorithms in training ViT model on CIFAR10 dataset. We plot the average stepsize $\frac{\gamma}{(\mathbf{D}_k)_{(j)}}$ (for Adam), $\frac{\tau_k}{(\mathbf{D}_k)_{(j)}}$ (for Momo-Adam), 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).

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

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

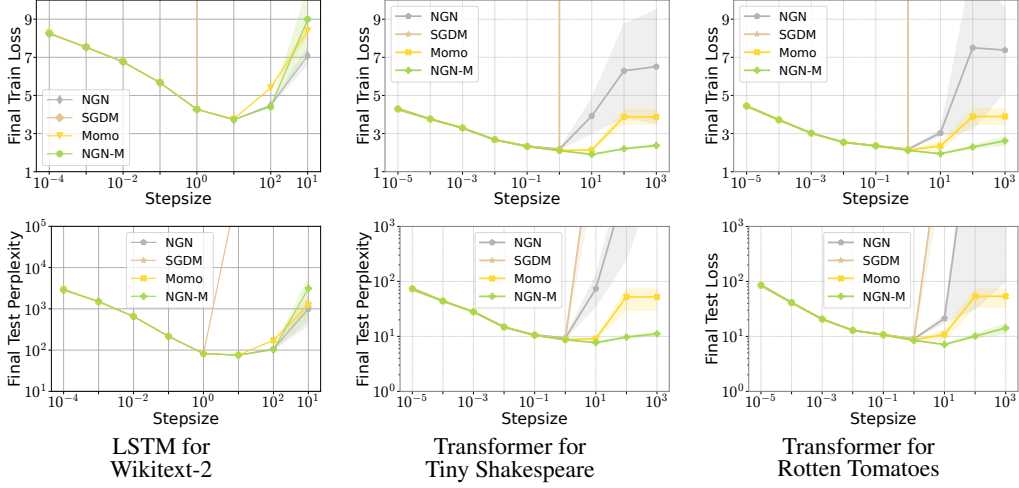


Figure J.8: Stability performance of algorithms supporting momentum and diagonal step-size varying step-size hyperparameter (c for NGN-M and NGN, α_0 for Momo, and step-size for SGDM). We observe that NGN-M achieves the training loss close to the best possible for a wider range of the step-size hyperparameter.

1430 in the beginning: the effective step-size increases at the beginning of the training. Then we observe
 1431 the main difference between the two algorithms above: effective step-size of Momo for sufficiently
 1432 large step-size hyperparameter is not adaptive within some part of the training, it always hits the
 1433 upper bound. Consequently, during that part of the training Momo reduces to SGDM. In contrast, the
 1434 effective step-size of NGN-M is always adaptive: it gradually decreases after a short warm-up. This
 1435 trend is similar to the state-of-the-art learning rate schedulers used in practice. Similar observations
 1436 can be made in comparison of NGN-MDv1 and Momo-Adam.

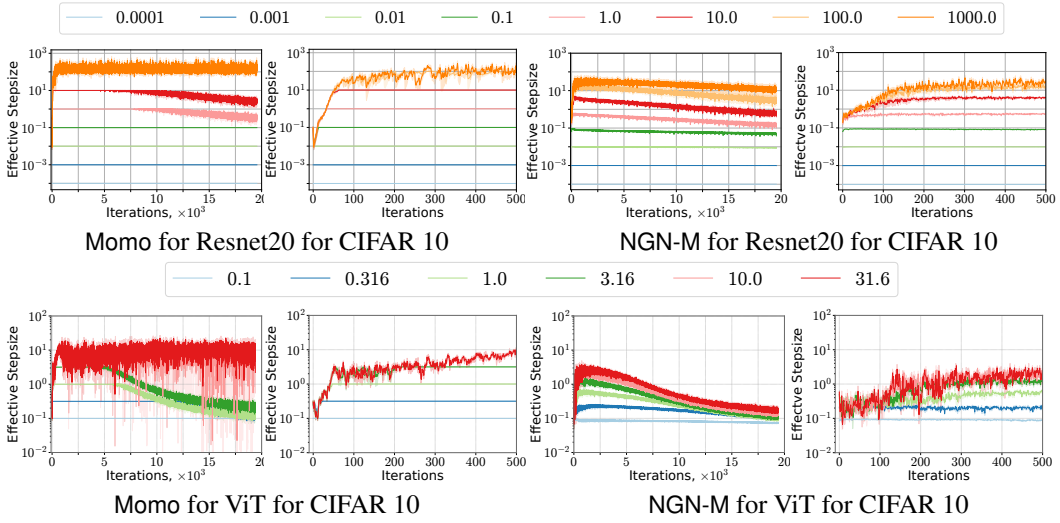


Figure J.9: The step-size of Momo and NGN-M during the training. We demonstrate the step-sizes τ_k for Momo and γ_k for NGN-M varying step-size parameters α_0 for Momo and c for NGN-M.

1437 J.10 Effective Updates in Training Language Models

1438 In this section, we demonstrate the magnitude of updates when training 160M language model with
 1439 Adam and NGN-MDv1 and varying the step-size hyperparameter across different layers of the model:
 1440 see the results in Figures J.11 to J.13. We demonstrate that NGN-MDv1 is a more conservative

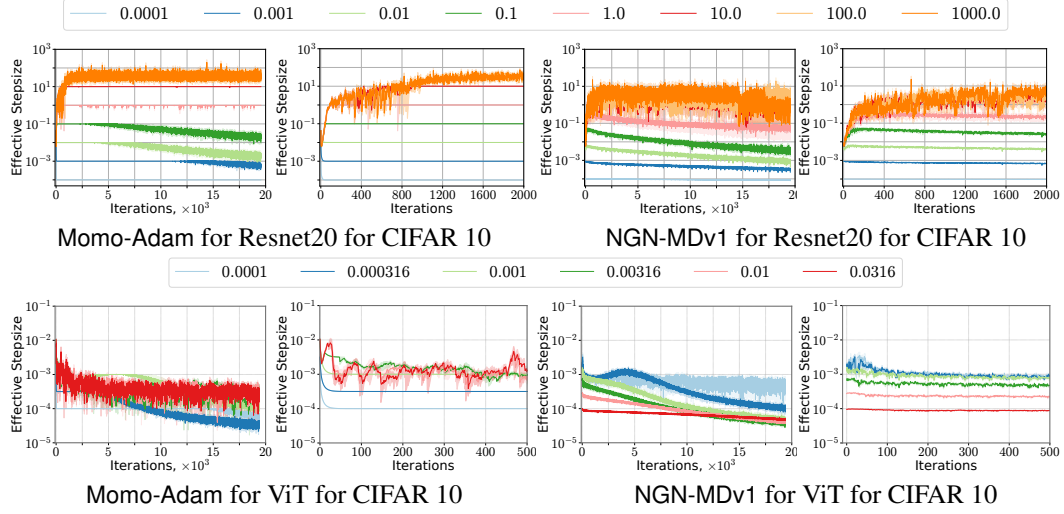


Figure J.10: The step-size of Momo-Adam and NGN-MDv1 during the training. We demonstrate the step-sizes τ_k for Momo-Adam and γ_k for NGN-MDv1 varying step-size parameters α_0 for Momo and c for NGN-MDv1.

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.

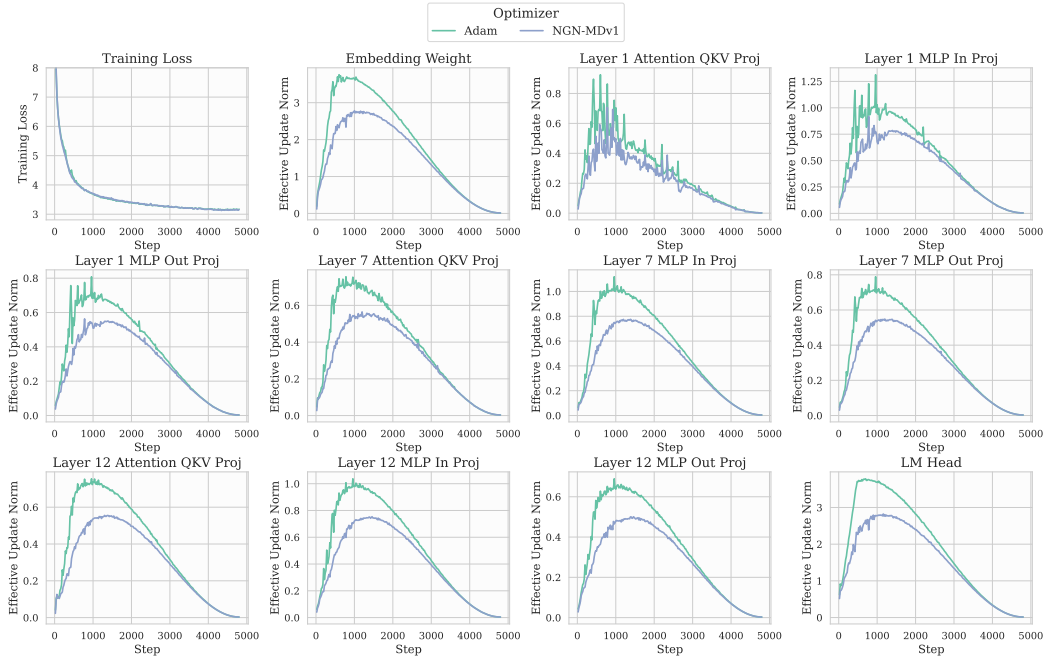


Figure J.11: Magnitude of updates when training 160M language model with Adam and NGN-MDv1 and step-size hyperparameter 0.003.

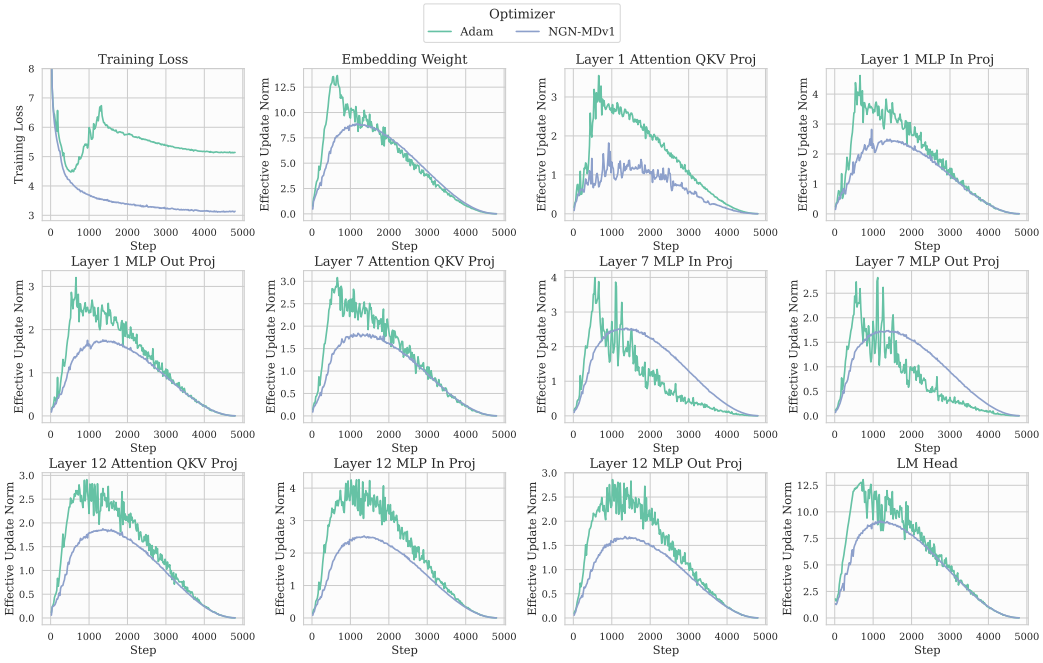


Figure J.12: Magnitude of updates when training 160M language model with Adam and NGN-MDv1 and step-size hyperparameter 0.01.

1446 **J.11 Training Dynamics in Training Language Models**

1447 Now we report the training dynamics in the training language across all tested sizes.

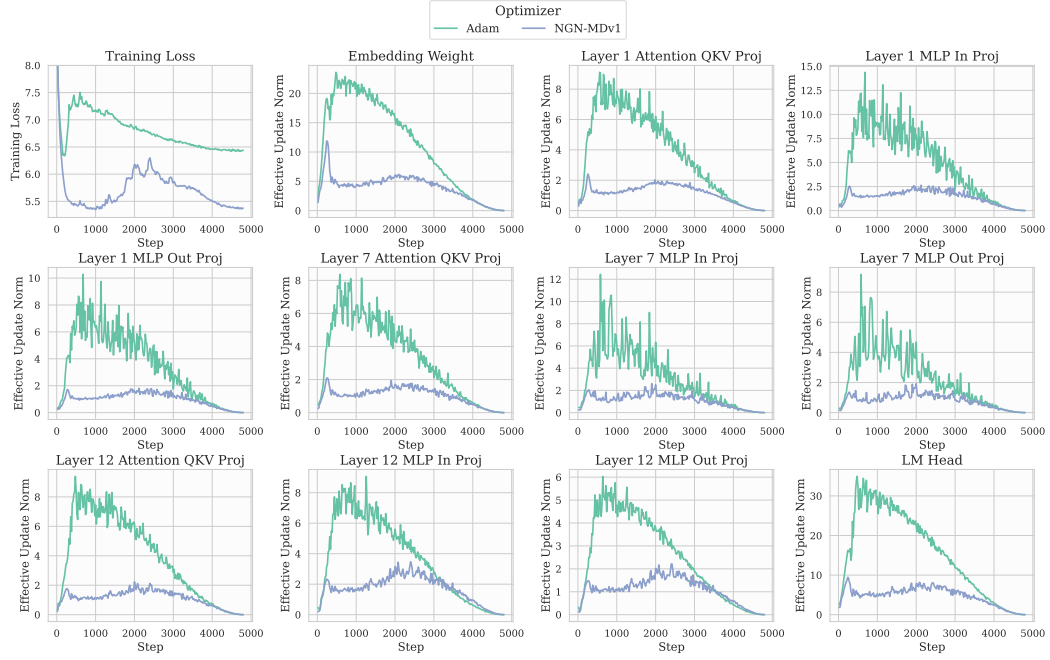


Figure J.13: Magnitude of updates when training 160M language model with Adam and NGN-MDv1 and step-size hyperparameter 0.03.

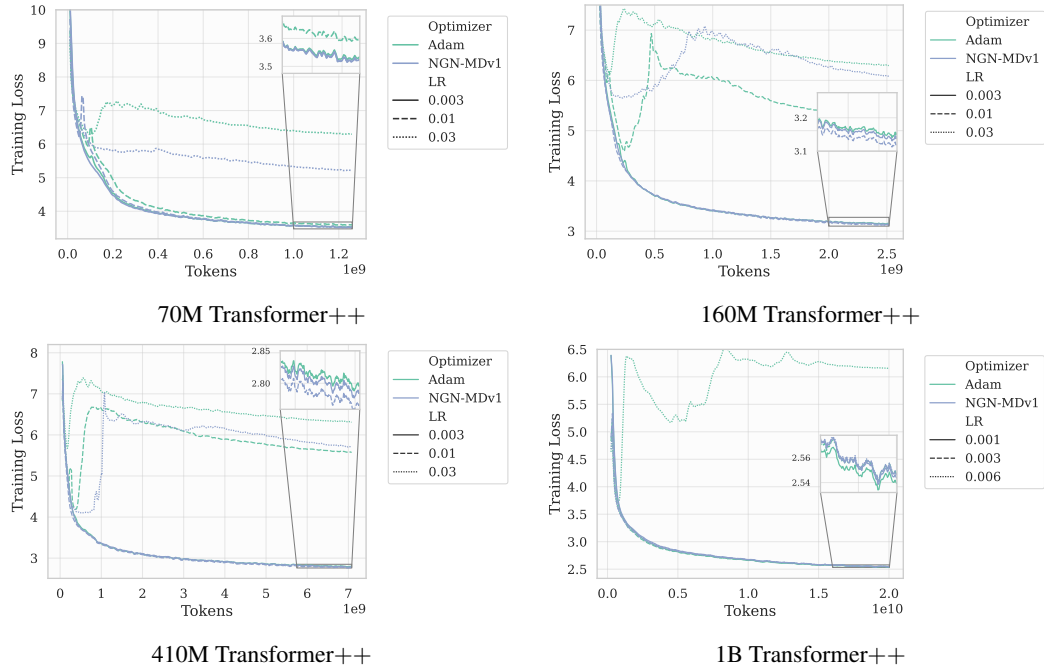


Figure J.14: Training dynamics when training language model at different sizes.