
Milo: A Magnitude-Invariant Learning Optimizer with Group-Wise Gradient Normalization

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 We introduce *Milo* (Magnitude-Invariant Learning Optimizer), a stochastic opti-
2 mizer that partitions gradients into network-wise or layer-wise groups and nor-
3 malizes each group to decouple update magnitudes from raw gradient norms. By
4 preventing overly aggressive steps in high-gradient regions and amplifying small
5 gradients, Milo produces smoother loss trajectories, more stable training, and
6 better generalization, while dramatically reducing sensitivity to the learning rate.
7 An optional scale-aware blending step preserves useful raw-gradient information
8 alongside invariance. We prove sublinear convergence in both smooth convex and
9 non-convex settings. Empirically, on synthetic loss landscapes, logistic regression,
10 multilayer perceptrons, ResNet-style CNNs, reinforcement-learning benchmarks,
11 and large-language-model pertaining, Milo matches or exceeds modern optimizers
12 in convergence speed and stability, cuts learning-rate tuning effort by up to 100 %,
13 and incurs only minimal runtime overhead. Our code is available at <URL>.

14 **1 Introduction**

15 Optimization is central to any decision-making problem, particularly modern-day machine learning
16 (ML) algorithms. Stochastic gradient descent (SGD) with momentum [Polyak, 1964, Sutskever et al.,
17 2013] and adaptive methods like AdamW [Loshchilov and Hutter, 2019, Kingma and Ba, 2015] have
18 become the most popular optimizers, with SGD performing exceedingly well in computer vision
19 [Goodfellow et al., 2016], and AdamW leading breakthroughs in Natural Language Processing (NLP),
20 large-scale transformers, and reinforcement learning [Vaswani et al., 2017, Haarnoja et al., 2018].
21 Yet despite their widespread prevalence, both optimizers exhibit weaknesses.

22 SGD’s effectiveness is highly contingent on initial learning rates and weight initializations, under-
23 mining its robustness across diverse tasks and architectures. Likewise, AdamW’s use of exponential
24 moving averages of squared gradients can produce unstable second-moment estimates, either vanish-
25 ing or exploding, which, while accelerating early convergence, often degrade ultimate generalization
26 performance. Due to these issues, in practical deep-learning workflows, substantial disparities in
27 gradient magnitudes across layers necessitate extensive per-layer hyperparameter searches, bespoke
28 (custom-made) learning-rate schedules, or reliance on gradient clipping and rescaling heuristics
29 (practical techniques).

30 These shortcomings lengthen hyperparameter tuning cycles, weaken training pipelines, and impede
31 reproducibility. For example, when training ResNet-50 on ImageNet, gradient norms in early
32 convolutional layers may be an order of magnitude smaller than those in the final classifier, causing
33 negligible updates in the former while the latter shifts erratically unless learning rates are calibrated
34 at each layer. In multi-task architectures, excessively large gradients from one branch can dominate
35 shared backbone updates, stalling progress on other tasks entirely.

36 **Contributions.** We aim to develop an optimizer that seamlessly handles gradient-scale heterogeneity
 37 (when gradient sizes vary, often by orders of magnitude, across different layers or parameters),
 38 combining the stability and generalization of SGD with the adaptability of AdamW, yet requiring
 39 minimal tuning.

40 To achieve this, we propose *Milo*, a magnitude-invariant learning optimizer that partitions gradients
 41 into network-wise (*Milo*) or layer-wise (*Milo*_{LW}) groups and normalizes each group before any
 42 momentum or adaptive scaling, in contrast to methods like NovoGrad [Ginsburg et al., 2019] which
 43 normalize moment estimates only after accumulation. We further introduce optional scale-aware
 44 blending to preserve useful raw-gradient information alongside invariance. Theoretically, we establish
 45 rigorous convergence guarantees in convex settings, with an $O(1/\sqrt{T})$ regret bound, and in non-
 46 convex settings, guarantees $\min_{t < T} \|\nabla F(w_t)\|^2 = O(T^{-1/2})$, meaning sublinear convergence to
 47 a first-order stationary point. Empirically, we validate *Milo* on six diverse tasks (synthetic landscapes,
 48 logistic regression, multilayer perceptrons (MLPs), Residual Neural Network with 18 layers (ResNet-
 49 18), reinforcement learning (RL), and large language model (LLM) pre-training), each averaged over
 50 five seeds with ablations on grouping size and blending. To ensure reproducibility, we provide the
 51 repository containing our testing, self-contained pseudocode, a comprehensive nomenclature, and
 52 details for hyperparameters, data-splits, and hardware. Across a wide range of benchmarks, *Milo*
 53 matches or exceeds the performance of modern optimizers, reduces learning-rate tuning effort by up
 54 to 100%, and adding minimal runtime overhead, making it a practical plug-and-play optimizer for
 55 modern deep-learning workflows.

Table 1: Averaged Train Loss, Test F1 Score, and Runtime by Optimizer and Model

Optimizer	Logistic			MLP			ResNet-18		
	Loss	F1	Time (s)	Loss	F1	Time (s)	Loss	F1	Time (s)
Milo	0.271	0.913	32.5	0.004	0.973	34.4	0.637	0.495	148.5
<i>Milo</i> _{LW}	0.226	0.919	32.7	0.020	0.969	36.1	1.419	0.417	182.8
SGD	0.277	0.915	29.7	0.073	0.969	30.3	0.572	0.448	130.6
AdamW	0.248	0.916	30.0	0.015	0.975	30.9	0.343	0.416	131.1
AdaGrad	0.468	0.879	29.8	0.196	0.942	30.7	0.468	0.433	130.7
NovoGrad	0.281	0.914	32.2	0.042	0.973	34.2	0.989	0.487	172.3

56 2 Related Work

57 **SGD and Momentum-Based Methods.** SGD [Robbins and Monro, 1951] iteratively updates
 58 parameters using noisy mini-batch gradients. While computationally efficient, it is sensitive to
 59 learning rates, can stall in saddle points, and may get trapped in local minima, although it often
 60 generalizes well [Ruder, 2016, Keskar et al., 2017]. In real-world training, this can manifest as
 61 extremely slow progress on deep models—e.g., early convolutional layers in ResNets barely update
 62 unless one carefully hand-tunes per-layer learning rates. To counter its slow convergence, momentum
 63 methods were introduced. These incorporate an exponential moving average of past gradients to
 64 smooth updates and accelerate convergence. Nesterov’s Accelerated Gradient (NAG) improves
 65 this further by computing the gradient at a look-ahead position [Nesterov, 1983]. Techniques like
 66 gradient normalization and centralization also help improve stability and convergence by recentring
 67 or scaling gradients across layers, but they still require careful hyperparameter choices to avoid under-
 68 or over-shooting [Yong et al., 2020, Xu et al., 2019].

69 **Adaptive Optimizers.** Adaptive optimizers adjust the learning rate on a per-parameter basis, en-
 70 hancing training stability and convergence [Duchi et al., 2011]. Early methods such as AdaGrad and
 71 RMSProp assign individual learning rates using past gradients [Duchi et al., 2011, Tieleman and
 72 Hinton, 2012], though they can overly reduce the effective rate in later stages of training, leading
 73 to vanishing updates on rarely-seen features [Wilson et al., 2017]. Adam [Kingma and Ba, 2015]
 74 combines momentum with adaptive scaling (tracking both first and second moments), making it less
 75 sensitive to hyperparameter tuning and widely popular in NLP and RL. However, Adam’s independent
 76 per-coordinate normalization can cause certain parameters to “overshoot” or overfit, especially in
 77 convolutional or recurrent architectures where gradients exhibit strong correlations [Heusel et al.,
 78 2017]. AdamW [Loshchilov and Hutter, 2019] improves on this by decoupling weight decay, enhanc-

79 ing generalization [Hettinger et al., 2022], but still inherits Adam’s sensitivity to extreme gradient
 80 magnitudes. NovoGrad [Ginsburg et al., 2019] builds on these ideas by applying layer-wise gradient
 81 normalization to the moment estimates, yielding stable, balanced updates even in high-variance or
 82 large-batch settings. Yet in practice it can still require manual tuning of subgroup sizes to avoid
 83 under-attenuation of noisy gradients.

84 3 Algorithm

Algorithm 1 Milo with median/MAD fallback and per-layer LR multipliers

Require: $\{p_i\}, \{g_i\}$, $\text{layer}(\cdot)$, $\eta, \{\mu_\ell\}, M, S, \alpha, \varepsilon$

```

1: for each layer  $\ell$  do                                ▷ process one layer at a time
2:    $v \leftarrow [\text{vec}(g_i) \mid \text{layer}(p_i) = \ell]$           ▷ collect & flatten layer gradients
3:    $m \leftarrow \max(1, \lceil |v|/M \rceil)$ ,  $s \leftarrow \lceil |v|/m \rceil$       ▷ compute #subgroups  $m$  and size  $s$ 
4:   for  $j = 0 \dots m - 1$  do                          ▷ normalize subgroup  $j$ 
5:      $v_j \leftarrow v[j \cdot s : \min((j + 1)s, |v|)]$           ▷ extract segment
6:     if  $\max|v_j| > 10^4$  or  $\min|v_j| < 10^{-6}$  then      ▷ detect outliers
7:        $\mu_j \leftarrow \text{median}(v_j)$                          ▷ robust center
8:        $d_j \leftarrow \text{median}(|v_j - \mu_j|)$                   ▷ MAD measure
9:        $\sigma_j \leftarrow 1.4826 d_j + \varepsilon$                    ▷ scale std
10:    else
11:       $\mu_j \leftarrow \text{mean}(v_j)$ ,  $\sigma_j \leftarrow \sqrt{\text{mean}((v_j - \mu_j)^2)} + \varepsilon$   ▷ standard mean/std
12:    end if
13:     $v_j \leftarrow (1 - S) \frac{v_j - \mu_j}{\sigma_j} + S v_j$           ▷ blend raw & normalized
14:  end for
15:  write back normalized  $v$  into  $\{g_i\}$                       ▷ update gradients in place
16: end for
17: for each parameter  $p_i$  do                            ▷ apply per-layer-scaled update
18:    $\ell \leftarrow \text{layer}(p_i)$ 
19:    $\gamma_i \leftarrow \eta \times (\mu_\ell \text{ if set else } 1)$       ▷ lookup layer LR multiplier
20:    $p_i \leftarrow p_i - \gamma_i g_i$                            ▷ SGD-style update
21: end for

```

85 3.1 Defining Milo

86 The key novelty of Milo is to *group gradients before normalization* and then apply per-layer learning-
 87 rate multipliers. Instead of normalizing the entire gradient or each coordinate individually, Milo
 88 partitions gradients into meaningful groups, normalizes each group independently (with an automatic
 89 median/median absolute deviation(MAD) fallback when needed), and finally scales each layer’s
 90 update by a user-specified multiplier. This yields balanced, outlier-resilient updates with minimal
 91 tuning.

92 Optionally, gradients may first be clipped:

$$g \leftarrow \frac{g}{\max(1, \|g\|_2/\text{clip_norm})} \quad (\text{gradient clipping}). \quad (3.1)$$

93 Given a parameter tensor p with raw gradient $g = \nabla f(p)$:

94 1. **Grouping.**

- 95 • *Layer-Wise (if enabled):*

$$\mathcal{G}_\ell = \{g_i \mid \text{layer}(p_i) = \ell\} \quad (\text{grouping by layer}). \quad (3.2)$$

- 96 • *Network-Wise:*

$$G = \text{reshape}(g, (-1, s)), \quad s = \max(1, \lfloor \sqrt{N} \rfloor) \quad (\text{flatten and chunk globally}). \quad (3.3)$$

97 **2. Group Normalization with Robust Fallback.** For each group G_k , compute mean and
 98 variance:

$$\mu_k = \frac{1}{|G_k|} \sum_{i \in G_k} g_i, \quad \sigma_k = \sqrt{\frac{1}{|G_k|} \sum_{i \in G_k} (g_i - \mu_k)^2 + \epsilon}. \quad (3.4)$$

99 If $\max_i |g_i| > 10^4$ or $\sigma_k < 10^{-6}$, use MAD:

$$\mu_k \leftarrow \text{median}(G_k), \quad \sigma_k \leftarrow 1.4826 \times \text{median}(|G_k - \mu_k|) + \epsilon. \quad (3.5)$$

100 Then normalize:

$$\tilde{g}_i = \frac{g_i - \mu_k}{\sigma_k} \quad \forall i \in G_k. \quad (3.6)$$

101 **3. Scale-Aware Blending (optional).**

$$\tilde{g}_i \leftarrow (1 - \alpha) \tilde{g}_i + \alpha g_i, \quad \alpha \in [0, 1] \quad (\text{preserve raw gradient structure}). \quad (3.7)$$

102 **4. Weight Decay.**

$$\tilde{g}_i \leftarrow \tilde{g}_i + \lambda p_i \quad (\text{L2 regularization}). \quad (3.8)$$

103 **5. Momentum & Adaptive Scaling.**

$$m_t = \beta m_{t-1} + \tilde{g} \quad (\text{momentum accumulation}), \quad (3.9)$$

$$\hat{m}_t = \begin{cases} \frac{m_t}{\sqrt{\sum_{\tau=0}^t m_{\tau}^2} + \text{adaptive_eps}}, & \text{if adaptive} = \text{True}, \\ m_t, & \text{otherwise} \end{cases} \quad (\text{adaptive RMS scaling}). \quad (3.10)$$

105 **6. Per-Layer LR Multipliers & Update.** Let γ be the base LR and μ_{ℓ} the per-layer multiplier:

$$\gamma_i = \gamma \times (\mu_{\ell} \text{ if defined, else 1}) \quad (\text{scale per layer}). \quad (3.11)$$

106 Then update parameters:

$$p \leftarrow p - \gamma_i \hat{m}_t \quad (\text{SGD-style update}). \quad (3.12)$$

107 **3.2 Group Normalization Methods**

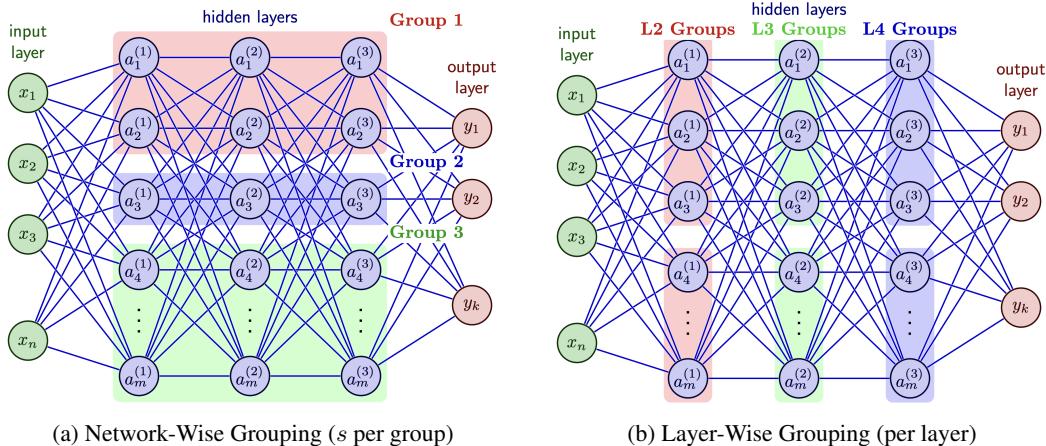


Figure 1: Comparison of Grouping Strategies

108 We propose two strategies to partition gradients for normalization:

109 **Network-Wise Grouping (Figure 1a).** Partition gradients into consecutive chunks of size s . Then
 110 apply mean/std normalization. Network-wise grouping, while reducing gradient variance through
 111 averaging, may disrupt intra-layer coordination when a single layer's parameters are spread across
 112 multiple groups, resulting in behavior closer to that of SGD.

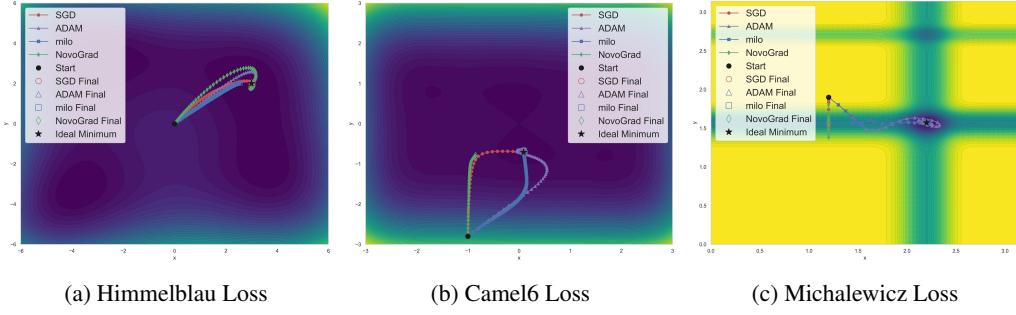


Figure 2: 2D contour loss maps for three benchmark functions

113 **Layer-Wise Grouping (Figure 1b).** Group all parameters within each layer, then normalize as above.
 114 Layer-wise grouping fixes the effective gradient norm (scaled by η) per layer, thereby reducing noise
 115 due to variance in gradient magnitudes. This normalization helps prevent any single layer from
 116 dominating the update dynamics, leading to smoother loss curves and promoting broader minima that
 117 may enhance generalization.

118 **Properties and Limitations** Milo’s per-group normalization yields the following properties. Vari-
 119 ance reduction results from the balanced updates by scaling down large-magnitude gradients and
 120 amplifying smaller ones. Magnitude invariance is an additional component from decoupling step
 121 size from raw gradient scale, in turn requiring less hyperparameter tuning. Scale-aware blending
 122 preserves raw gradient information based on the desired. Additionally, there is an architectural
 123 agnosticism which works across heterogeneous layer types with minimal sensitivity. The optimizer is
 124 also more robust, automatically switching to a median/MAD-based normalization when encountering
 125 extreme gradient values. Lastly, Milo has demonstrated smoother and more direct convergence when
 126 compared to other optimizers. Some negative properties of Milo include dense-update sensitivity and
 127 a consistency versus adaptivity trade-off. There is an inherent grouping strategy dependency, with
 128 performance critically depending on how gradients are partitioned. Additionally, the extra grouping
 129 normalization steps can add runtime or memory overhead.

130 **Algorithm Intuition** To gain intuition, we compare Milo with SGD, Adam, and NovoGrad on three-
 131 dimensional loss landscapes. On Himmelblau [Himmelblau, 1972] (Figure 2a), all methods reach the
 132 global minimum, but Milo and SGD follow smooth, direct paths while Adam and NovoGrad oscillate.
 133 In the Camel6 landscape [Molga, 2005] (Figure 2b), Milo’s uniform step size produces steady
 134 descent, avoiding the initial overshooting seen in the others. Finally, on Michalewicz [Michalewicz,
 135 1996] (Figure 2c), only Milo and Adam escape the shallow basin to reach the deeper minimum;
 136 SGD and NovoGrad remain trapped. These observations highlight that Milo’s fixed step size and
 137 group-wise normalization reduce variance and oscillations, striking a balance between stable progress
 138 and adaptive exploration. Overall, Milo’s group-wise normalization was shown to blend SGD’s stable,
 139 direct convergence with Adam’s adaptability. It yields low-variance trajectories that escape shallow
 140 valleys and avoid the oscillations typical of fully adaptive methods.¹

141 **Efficient Implementation** A large focus during the development of Milo was to have an efficient
 142 implementation due to the inherent cost of partitioning and normalizing gradients. Our implemen-
 143 tation greatly minimizes the overhead that comes with this approach, implementing the following.
 144 Firstly, we minimize extra memory allocations by copying gradients into a pre-allocated buffer (stored
 145 in `_cached_buffers`). We also perform vectorized and batched computations of group mean and
 146 standard deviations to ensure a minimal memory footprint. As layer-wise is more computationally
 147 expensive, we use cached layer mapping to avoid repeated computations and overhead, generating a
 148 mapping of parameters to layers via the `_organize_layer_groups` function.

149 **Difference from NovoGrad** Milo partitions raw gradients into network or layer-wise groups and
 150 normalizes each group before any momentum or adaptive scaling—decoupling update size from
 151 gradient magnitude and offering optional scale-aware blending α . NovoGrad instead normalizes

¹See Appendix D.5 for precise equations, more in-depth analysis, and three-dimensional plots.

152 per-layer moment estimates after accumulation, so its step sizes remain tied to those norms and
 153 lack explicit blending. Consequently, Milo is plug-and-play with minimal tuning, while NovoGrad
 154 requires careful adjustment of learning rate, decay rates, and subgroup sizes.

155 4 Convergence Guarantees for Milo

156 In this section, we establish theoretical convergence guarantees for Milo under a non-convex setting.
 157 We show that Milo benefits from variance reduction via layer-wise grouping. This variance reduction
 158 not only improves the effective gradient norm, yielding lower regret constants, but also has important
 159 implications for generalization across a variety of tasks. Due to space constraints, the detailed convex
 160 convergence proofs are deferred to Appendix B. Here, we focus on the non-convex setting.²

161 4.1 Convergence in Non-Convex Settings

162 Under the same normalization and L -smoothness assumptions, one shows:

163 *Proof Sketch of Theorem C.* Under Assumptions C.1, C.2, and the variance contraction/moment
 164 conditions (Assumptions B.3, B.4), we outline the main steps.

165 **1. Descent Lemma.** By L -smoothness of F , for $w_{t+1} = w_t - \eta \tilde{g}_t$:

$$F(w_{t+1}) \leq F(w_t) - \eta \nabla F(w_t), \tilde{g}_t + \frac{L\eta^2}{2} \|\tilde{g}_t\|^2. \quad (4.1)$$

166 **2. Alignment Bound.** Decompose $\tilde{g}_t = \nabla F(w_t) + \delta_t$. From Lemma B.2,

$$\nabla F(w_t), \tilde{g}_t \geq (1 - \varepsilon) \|\nabla F(w_t)\|^2.$$

167 **3. Norm Control.** Using variance contraction (Lemma B.4),

$$\|\tilde{g}_t\|^2 \leq \rho n.$$

168 **4. One-Step Inequality.** Combine the above into the descent:

$$F(w_{t+1}) \leq F(w_t) - \eta(1 - \varepsilon) \|\nabla F(w_t)\|^2 + \frac{L\eta^2}{2} \rho n. \quad (4.2)$$

169 Rearrange:

$$\eta(1 - \varepsilon) \|\nabla F(w_t)\|^2 \leq F(w_t) - F(w_{t+1}) + \frac{L\eta^2}{2} \rho n.$$

170 **5. Telescoping Sum.** Summing for $t = 0 \dots T - 1$ and using $F(w_T) \geq F^*$:

$$\sum_{t=0}^{T-1} \|\nabla F(w_t)\|^2 \leq \frac{F(w_0) - F^*}{\eta(1 - \varepsilon)} + \frac{L\eta T}{2(1 - \varepsilon)} \rho n. \quad (4.3)$$

171 **6. Rate via Step-Size.** Choose $\eta = \min\{1/L, c/\sqrt{T}\}$ for constant $c > 0$. Then both terms scale
 172 as $O(T^{-1/2})$, yielding

$$\min_{t < T} \|\nabla F(w_t)\|^2 = O(T^{-1/2}). \quad (4.4)$$

173 \square

174 By normalizing groups of gradient components, Milo achieves variance contraction, which in turn
 175 improves convergence rates compared to standard SGD in both convex and non-convex settings.
 176 Additionally, the variance contraction also stabilizes the optimization trajectory. This stability can
 177 serve as an implicit regularizer, biasing the algorithm toward flatter minima, a trait empirically linked
 178 to improved generalization performance [Gower et al., 2020]. Additionally, a discussion of key
 179 technical aspects in these theorems is given in ??.

²For full convergence proofs, see Appendix B for convex and Appendix C for non-convex. Here we have precise equations, review key convexity properties, derive regret bounds for standard SGD, and then show that Milo benefits from variance reduction via layer-wise grouping.

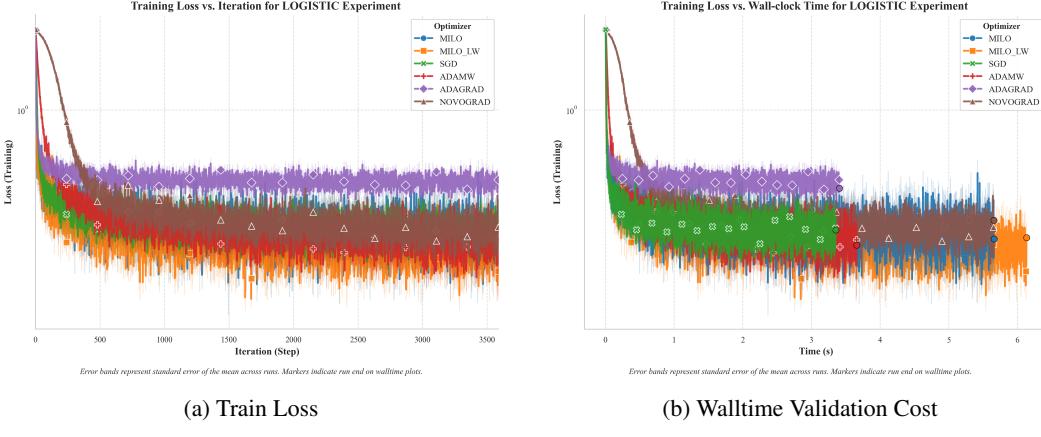


Figure 3: Comparison of Logistic Regression training loss by iterations and wall time.

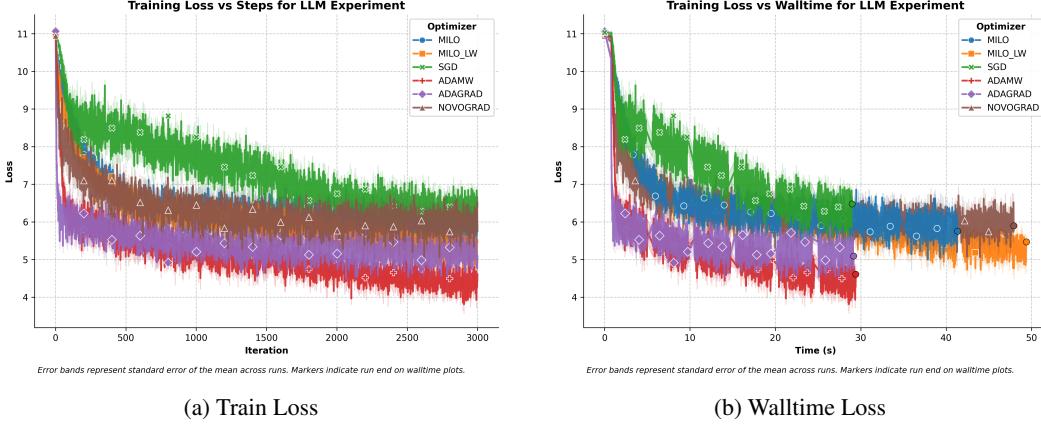


Figure 4: Comparison of LLM performance by training loss and wall time.

180 5 Experiments

181 To evaluate Milo empirically, we compare it with SGD with momentum, AdamW, AdaGrad, and
 182 NovoGrad across six experiments to cover common machine learning models. These experiments
 183 include finding the global minima for loss maps, logistic regression, MLPs, convolutional neural
 184 networks (CNNs) [LeCun et al., 1998] via ResNet-18 [He et al., 2016], RL in a three-dimensional
 185 space, and LLM pre-training. The models were trained in 5 separate instances for each empirical
 186 experiment, and the results of these trainings were averaged. The same parameter initialization
 187 is used when comparing the optimizers, and an extensive hyperparameter search and tuning were
 188 performed individually for each optimizer, as shown in Appendix D.6.³ All Adam variants use default
 189 parameters $\beta_1 = 0.9$, $\beta_2 = 0.999$, $\epsilon = 10^{-8}$. This experimentation demonstrates that Milo can solve
 190 deep-learning problems with a high loss reduction and convergence rate.

191 5.1 Experimentation Analysis

192 Our experimentation found that across diverse tasks, Milo and its layer-wise variant (Milo_{LW}) consis-
 193 tently achieved improved or on-par convergence and enhanced stability over traditional optimizers.
 194 It can be observed that the normalized gradient updates reduce the variance and noise, while also
 195 ensuring invariance to raw gradient scales. These overall contribute to its consistent progress in both
 196 convex and non-convex settings. However, its walltime is consistently worse or on par with the other
 197 evaluated optimizers.

³See Appendix D for details of all hyperparameters and their tuning, data splits, hardware used, compute resources, in-depth statistical significance tests, experiment accuracy results, ablation study, and link to test code.

Average Agent Trajectories by Optimizer

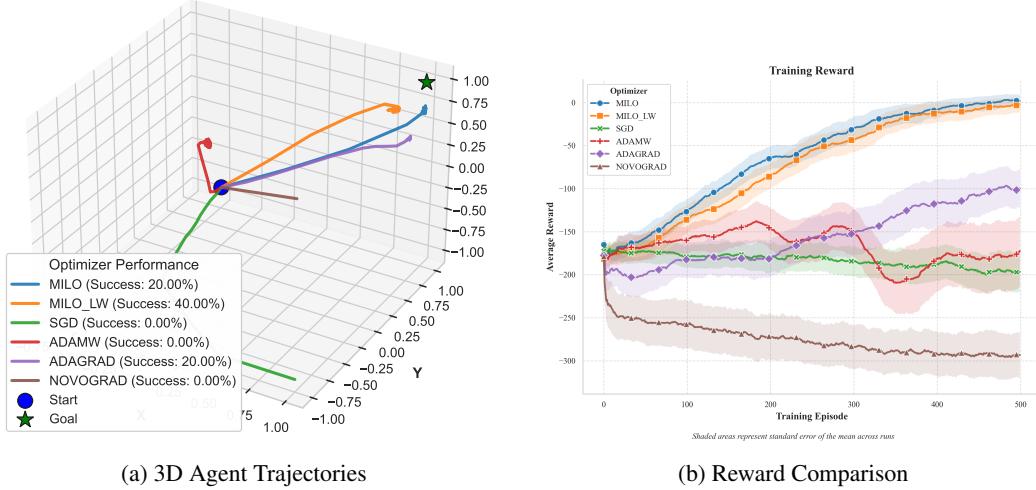


Figure 5: Comparison of RL results in terms of agent trajectories and reward performance.

- 198 In the logistic regression and MLPs experiments, shown in Appendix D.7 and Appendix D.8, Milo_{LW}
 199 achieves faster loss reduction along with higher F1 accuracy in convex and non-convex settings with
 200 these less complex problems (Figure 3a and Figure 3b). This illustrates the benefits of the layer-wise
 201 group-normalized updates in these less complex environments, leading to a more rapid loss decrease.
 202 The ResNet18 experiment (in Appendix D.9) contrasts with the less complex experiments as Milo
 203 outperforms Milo_{LW}. Milo demonstrates robust performance with an almost on-par convergence rate
 204 with its competitors, while also demonstrating a high validation accuracy (Table 1). This suggests
 205 that while Milo_{LW} is effective for capturing layer-specific dynamics, it can underperform in scenarios
 206 where global gradient invariance is more advantageous (e.g., CNNs [Ye et al., 2021]).
 207 In the LLM experiment (in Appendix D.11), Milo_{LW} is one of the better performing optimizers,
 208 demonstrating robust performance (Figure 4a and Figure 4b). However, it falls close behind AdamW
 209 and AdaGrad. This suggests that global grouping can wash out useful layer disparities in particular
 210 complex instances such as the LLM pre-training here.
 211 In Appendix D.10, we evaluate the optimizers in a three-dimensional RL scenario (Figure 5a and
 212 Figure 5b). For this task, there were two variations, featuring the agents reaching a consistently set
 213 point and a randomly generated point. In the set goal task, most optimizers had a near 100% accuracy
 214 in reaching the goal. Whereas in the randomized setup, Milo_{LW} achieves the highest success rate
 215 by the largest margin, with Milo, and AdaGrad following behind, with half of its success rate. This
 216 demonstrates the better performance of Milo_{LW} in noisy or ill-conditioned tasks, with it also best
 217 maximizing the reward for these scenarios as well.
 218 **Scalability and Runtime.** This study also evaluated the scalability and efficiency of Milo and
 219 Milo_{LW}, with it being shown that in a large variety of scenarios, their wall-clock time, memory
 220 usage, and computational overhead often come at the cost of their increase in performance. Though
 221 Milo was observed to maintain relatively on-par runtimes with other optimizers, Milo_{LW} commonly
 222 sacrificed its runtime for better results in the more complex models, such as the LLM pre-training.
 223 Future work will involve addressing these runtime issues and further expanding the initial solutions
 224 in Section 3.2. A more in-depth analysis of the runtime metrics is provided in Appendix D.
 225 **Hyperparameter Tuning Behavior.** Despite aggressive hyperparameter tuning, Milo and Milo_{LW}
 226 proved to be largely invariant to changes in their hyperparameters. Due to this, the default values for
 227 Milo and Milo_{LW} were employed for each experiment, with no tuning. This suggests an invariance
 228 to parameter changes, with the normalization acting as a self-regulator for the optimizer, giving it
 229 a "plug and play" usage that many other optimizers do not provide. In contrast, the other methods
 230 required careful tuning to reach the demonstrated accuracies, with their default implementation
 231 achieving much lower accuracies. We provide further exploration of this idea in Appendix D.12.

232 **5.2 Theory to Empirical Results**

233 In convex settings, experiments on logistic regression have demonstrated a monotonic decrease in
234 training loss consistent with our descent inequality. The normalized updates enforce a stable step size
235 regardless of raw gradient magnitude variations. Additionally, our theoretical analysis (Appendix B)
236 shows that controlling the effective gradient norm (via the contraction factor ρ) and the alignment
237 error (ε) reduces the regret from $R_{\text{SGD}}(T) \leq D G \sqrt{T}$ to $R_{\text{Milo}}(T) \leq D \tilde{G} \sqrt{T}$, with $\tilde{G} < G$.
238 Empirical results (??) confirm that both Milo and its layer-wise variant (Milo_{LW}) consistently achieve
239 lower training loss and higher accuracy on MNIST.
240 In non-convex settings, stable convergence trajectories were given for complex models like multilayer
241 neural networks and CNNs, as our non-convex convergence theorem (Appendix C) guarantees that
242 $\min_{t=1,\dots,T} \mathbb{E} [\|\nabla F(w_t)\|^2] = O\left(\frac{1}{\sqrt{T}}\right)$, implying reduced gradient noise. Empirically, Figure 4a
243 and Figure 4b show that Milo-based updates yield smoother loss curves and stable convergence
244 despite the presence of multiple local minima and saddle points. Furthermore, the group normalization
245 reduces gradient variance (as quantified by ρ and controlled via ϵ and σ_{\min}), which is reflected in
246 lower oscillations and tighter convergence in the loss curves.

247 **6 Limitations and Future Work**

248 **Limitations.** While Milo demonstrates strong performance across tasks, several important questions
249 remain. First, its behavior in large-scale fine-tuning regimes (e.g. LLM fine-tuning) has yet to be
250 explored. Second, in densely parameterized modules, such as attention layers, per-group normaliza-
251 tion may inadvertently suppress critical gradient signals. Last, by enforcing nearly constant step
252 magnitudes, Milo trades off adaptivity for consistency; this could lead to overly conservative updates
253 in rapidly changing or highly nonstationary loss landscapes.

254 **Extensions and Variants.** We believe our algorithm could be extended to cover and further evaluate
255 these limitations in future work. Such extensions would include addressing the above limitations,
256 exploring new versions such as Milo with decoupled weight decay, and further theoretical analysis
257 to compare the similarities of base Milo with SGD and Milo_{LW} with Adam. We plan to expand our
258 implementation to cover those mentioned in future work.

259 **7 Conclusion**

260 This study introduces a new, high-performing optimizer applicable to diverse ML training scenarios.
261 The major impact of our work is the increased stability, loss reduction, and reduced sensitivity to
262 hyperparameters. These make it practical for practitioners looking for a "plug and play" optimizer
263 that can retain or exceed the performance of other optimizers. Milo combines the stability and
264 direct update dynamics of SGD with the variance-reducing effects of gradient normalization, making
265 it effective across vision tasks, varying model scales, and noisy or ill-conditioned optimization
266 settings. Furthermore, Milo_{LW} offers enhanced layer-specific adaptation, excelling in settings where
267 fine-grained adjustments are required due to heterogeneous gradient landscapes. It is to be noted that
268 despite the improved performance, Milo comes at the cost of increased runtime, especially when
269 considering the layer-wise version. Overall, the Magnitude-Invariant Learning Optimizer (Milo and
270 Milo_{LW}) provides a robust, theoretically sound, and practically effective alternative to conventional
271 optimizers.

272 **References**

- 273 Akiba, T., Sano, S., Yanase, T., Ohta, T., and Koyama, M. (2019). Optuna: A next-generation hyperparameter
274 optimization framework. In *Proceedings of the 25th ACM SIGKDD International Conference on Knowledge
275 Discovery and Data Mining*, pages 2623–2631.
- 276 Bottou, L., Curtis, F. E., and Nocedal, J. (2018). Optimization methods for large-scale machine learning. *SIAM
277 Review*, 60(2):223–311. Used for Assumption 1.2 (Bounded Gradient Norm).
- 278 Brockman, G., Cheung, V., Pettersson, L., Schneider, J., Schulman, J., Tang, F., and Zaremba, W. (2016). Openai
279 gym. arXiv:1606.01540.
- 280 Crockford, D. (2006). The application/json media type for javascript object notation (json). RFC 4627.
- 281 Duchi, J., Hazan, E., and Singer, Y. (2011). Adaptive subgradient methods for online learning and stochastic
282 optimization. *Journal of Machine Learning Research*, 12:2121–2159. Introduces AdaGrad and provides a
283 detailed regret analysis of adaptive gradient methods, closely related to Adam’s approach.
- 284 FareedKhan-dev (2025). FareedKhan-dev/train-lm-from-scratch: A straightforward method to train large
285 language models from scratch. GitHub repository. Accessed: 2025-05-14.
- 286 Gao, L., Biderman, S., Black, S., Golding, C., Hoppe, T., Foster, J. P., He, H., Thite, A., Nabeshima, N.,
287 Tow, B., et al. (2020). The pile: An 800gb dataset of diverse text for language modeling. *arXiv preprint
288 arXiv:2101.00027*.
- 289 Ginsburg, B., Castonguay, P., Hrinchuk, O., Kuchaiev, O., Lavrukhin, V., Leary, R., Li, J., Nguyen, H., Zhang,
290 Y., and Cohen, J. M. (2019). Stochastic gradient methods with layer-wise adaptive moments for training of
291 deep networks. *Preprint at https://arxiv.org/abs/1905.11286*.
- 292 Goodfellow, I., Bengio, Y., and Courville, A. (2016). *Deep Learning*. MIT Press.
- 293 Gower, R. M., Schmidt, M., Bach, F., and Richtárik, P. (2020). Variance-reduced methods for machine learning.
294 *Proceedings of the IEEE*, 108(11):1968–1983. Shows empirical affect of variance reduction.
- 295 Haarnoja, T., Zhou, A., Abbeel, P., and Levine, S. (2018). Soft actor-critic: Off-policy maximum entropy deep
296 reinforcement learning with a stochastic actor. In *International Conference on Machine Learning (ICML)*,
297 pages 1861–1870.
- 298 He, K., Zhang, X., Ren, S., and Sun, J. (2016). Deep residual learning for image recognition. In *Proceedings of
299 the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 770–778.
- 300 Hettinger, L., Bernstein, J., et al. (2022). Recent advances in optimizers for deep learning. *arXiv preprint
301 arXiv:2207.01716*.
- 302 Heusel, M., Ramsauer, H., Unterthiner, T., Nessler, B., and Hochreiter, S. (2017). Gans trained by a two
303 time-scale update rule converge to a local nash equilibrium. In *Advances in Neural Information Processing
304 Systems*, volume 30.
- 305 Himmelblau, D. M. (1972). *Applied Nonlinear Programming*. McGraw-Hill.
- 306 Hoeffding, W. (1963). Probability inequalities for sums of bounded random variables. *Journal of the American
307 Statistical Association*, 58(301):13–30. Used for Lemma 1.1 (Concentration of Empirical Means).
- 308 Hunter, J. D. (2007). Matplotlib: A 2d graphics environment. *Computing in Science & Engineering*, 9(3):90–95.
- 309 Ioffe, S. and Szegedy, C. (2015). Batch normalization: Accelerating deep network training by reducing internal
310 covariate shift. In *Proceedings of the 32nd International Conference on Machine Learning (ICML)*, pages
311 448–456. Used for Lemma 1.4 (Bounded Variance Contraction).
- 312 Keskar, N. S., Nocedal, J., et al. (2017). On large-batch training for deep learning: Generalization gap and sharp
313 minima. In *International Conference on Learning Representations*.
- 314 Kingma, D. P. and Ba, J. (2015). Adam: A method for stochastic optimization. In *International Conference
315 on Learning Representations (ICLR)*. Introduces Adam, which computes adaptive learning rates for each
316 parameter by combining momentum and adaptive scaling of gradients. This approach reduces variance and
317 accelerates convergence compared to plain SGD.
- 318 Krizhevsky, A. and Hinton, G. (2009). Learning multiple layers of features from tiny images. Technical report,
319 University of Toronto.

- 320 LeCun, Y., Bottou, L., Bengio, Y., and Haffner, P. (1998). Gradient-based learning applied to document
 321 recognition. *Proceedings of the IEEE*, 86(11):2278–2324. The MNIST dataset consists of 70,000 grayscale
 322 images of handwritten digits (0–9) and is commonly used for training and testing in the field of machine
 323 learning and computer vision.
- 324 lonePatient (2019). lonePatient/NovoGrad-pytorch: Pytorch implementation of the novograd optimizer. GitHub
 325 repository. Accessed: 2025-05-14.
- 326 Loshchilov, I. and Hutter, F. (2019). Decoupled weight decay regularization. In *International Conference on
 327 Learning Representations*.
- 328 McKinney, W. (2010). Data structures for statistical computing in python. In *Proceedings of the 9th Python in
 329 Science Conference*, pages 51–56.
- 330 Michalewicz, Z. (1996). *Genetic Algorithms + Data Structures = Evolution Programs*. Springer, 3rd edition.
- 331 Molga, M. and Smutnicki, C. (2005). Test functions for optimization needs. Technical report, Faculty of
 332 Fundamentals of Technology, Wrocław University of Technology.
- 333 Nesterov, Y. (1983). A method for unconstrained convex minimization problem with the rate of convergence
 334 $o(1/k^2)$. *Soviet Mathematics Doklady*, 27(2):372–376. Used for Assumption 1.1 (Smoothness).
- 335 Paszke, A., Gross, S., Massa, F., Lerer, A., Bradbury, J., Chanan, G., Killeen, T., Lin, Z., Gimelshein, N., Antiga,
 336 L., Desmaison, A., Kopf, A., Yang, E., DeVito, Z., Raison, M., Tejani, A., Chilamkurthy, S., Steiner, B., Fang,
 337 L., Bai, J., and Chintala, S. (2019). Pytorch: An imperative style, high-performance deep learning library.
 338 In Wallach, H. M., Larochelle, H., Beygelzimer, A., d’Alché Buc, F., Fox, E. B., and Garnett, R., editors,
 339 *Advances in Neural Information Processing Systems 32*, pages 8024–8035. Curran Associates, Inc.
- 340 Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., Blondel, M., Prettenhofer, P.,
 341 Weiss, R., Dubourg, V., Vanderplas, J., Passos, A., Cournapeau, D., Brucher, M., Perrot, M., and Duchesnay,
 342 (2011). Scikit-learn: Machine learning in python. *Journal of Machine Learning Research*, 12:2825–2830.
- 343 Polyak, B. T. (1964). Some methods of speeding up the convergence of iteration methods. *USSR Computational
 344 Mathematics and Mathematical Physics*, 4(5):1–17.
- 345 Python Software Foundation (2023). *Python Standard Library: itertools*. Version 3.10.
- 346 PyTorch Contributors (2023). torchvision: image and video datasets, transforms, and models for pytorch.
 347 <https://github.com/pytorch/vision>. Version 0.21.0.
- 348 Robbins, H. and Monroe, S. (1951). A stochastic approximation method. *The Annals of Mathematical Statistics*,
 349 22(3):400–407.
- 350 Ruder, S. (2016). An overview of gradient descent optimization algorithms. *arXiv preprint arXiv:1609.04747*.
- 351 Sutskever, I., Martens, J., Dahl, G., and Hinton, G. (2013). On the importance of initialization and momentum in
 352 deep learning. In *International Conference on Machine Learning (ICML)*, pages 1139–1147.
- 353 Tieleman, T. and Hinton, G. (2012). Lecture 6.5—rmsprop: Divide the gradient by a running average of its
 354 recent magnitude. COURSERA: Neural Networks for Machine Learning.
- 355 van der Walt, S., Colbert, S. C., and Varoquaux, G. (2011). The numpy array: A structure for efficient numerical
 356 computation. *Computing in Science & Engineering*, 13(2):22–30.
- 357 Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., Kaiser, L., and Polosukhin, I.
 358 (2017). Attention is all you need. In *Advances in Neural Information Processing Systems (NeurIPS)*.
- 359 Vershynin, R. (2018). *High-Dimensional Probability: An Introduction with Applications in Data Science*.
 360 Cambridge University Press. Used for Assumption 1.4 (Moment Conditions).
- 361 Virtanen, P., Gommers, R., Oliphant, T. E., Haberland, M., Reddy, T., Cournapeau, D., Burovski, E., Peterson, P.,
 362 Weckesser, W., Bright, J., et al. (2020). Scipy 1.0: Fundamental algorithms for scientific computing in python.
 363 *Nature Methods*, 17:261–272.
- 364 Waskom, M. (2021). Seaborn: Statistical data visualization.
- 365 Wilson, A. C., Roelofs, R., Stern, M., Srebro, N., and Recht, B. (2017). The marginal value of adaptive gradient
 366 methods in machine learning. In *Advances in Neural Information Processing Systems*, volume 30.

- 367 Wolf, T., Debut, L., Sanh, V., Chaumond, J., Delangue, C., Moi, A., Cistac, P., Rault, T., Louf, R., Funtowicz, M.,
368 and Brew, J. (2020). Transformers: State-of-the-art natural language processing. In *Proceedings of the 2020*
369 *Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 38–45.
- 370 Wu, Y. and He, K. (2018). Group normalization. In *Proceedings of the European Conference on Computer*
371 *Vision (ECCV)*, pages 3–19. Used for Assumption 1.3 (Variance Contraction via Normalization) and Lemma
372 1.3 (Robustness for Small σ_j).
- 373 Xu, J., Sun, X., Zhang, Z., Zhao, G., and Lin, J. (2019). Understanding and improving layer normalization. In *Advances in Neural Information Processing Systems*, volume 32, pages 5665–5675. Curran Associates, Inc.
374 We reference this work to highlight limitations of per-layer normalization in handling heterogeneous gradient
375 scales, motivating our group-normalization approach in Milo.
- 376 Ye, C., Zhou, X., McKinney, T., Liu, Y., Zhou, Q., and Zhdanov, F. (2021). Exploiting invariance in training
377 deep neural networks. In *Proceedings of the Thirty-Eighth AAAI Conference on Artificial Intelligence*. This
378 research paper delves into the theoretical advantages of incorporating invariance into neural networks, a
379 common characteristic of real-world data.
- 380 Yong, Y., Huang, E., Hua, L., and Zhang, C. (2020). Gradient centralization: A new optimization technique
381 for deep neural networks. In *Advances in Neural Information Processing Systems (NeurIPS)*. Introduces
382 Gradient Centralization, an optimization technique that centralizes gradients to improve training stability and
383 convergence.

385 **NeurIPS Paper Checklist**

386 The checklist is designed to encourage best practices for responsible machine learning research, addressing
387 issues of reproducibility, transparency, research ethics, and societal impact. Do not remove the checklist: **The**
388 **papers not including the checklist will be desk rejected.** The checklist should follow the references and follow
389 the (optional) supplemental material. The checklist does NOT count towards the page limit.

390 Please read the checklist guidelines carefully for information on how to answer these questions. For each
391 question in the checklist:

- 392 • You should answer **[Yes]**, **[No]**, or **[NA]**.
393 • **[NA]** means either that the question is Not Applicable for that particular paper or the relevant
394 information is Not Available.
395 • Please provide a short (1–2 sentence) justification right after your answer (even for NA).

396 **The checklist answers are an integral part of your paper submission.** They are visible to the reviewers, area
397 chairs, senior area chairs, and ethics reviewers. You will be asked to also include it (after eventual revisions)
398 with the final version of your paper, and its final version will be published with the paper.

399 The reviewers of your paper will be asked to use the checklist as one of the factors in their evaluation. While
400 "**[Yes]**" is generally preferable to "**[No]**", it is perfectly acceptable to answer "**[No]**" provided a proper
401 justification is given (e.g., "error bars are not reported because it would be too computationally expensive" or
402 "we were unable to find the license for the dataset we used"). In general, answering "**[No]**" or "**[NA]**" is not
403 grounds for rejection. While the questions are phrased in a binary way, we acknowledge that the true answer is
404 often more nuanced, so please just use your best judgment and write a justification to elaborate. All supporting
405 evidence can appear either in the main paper or the supplemental material, provided in appendix. If you answer
406 **[Yes]** to a question, in the justification please point to the section(s) where related material for the question can
407 be found.

408 **IMPORTANT**, please:

- 409 • **Delete this instruction block, but keep the section heading "NeurIPS Paper Checklist",**
410 • **Keep the checklist subsection headings, questions/answers and guidelines below.**
411 • **Do not modify the questions and only use the provided macros for your answers.**

412 **1. Claims**

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

415 Answer: **[Yes]**

416 Justification: This paper provides extensive empirical and theoretical validation, which matches
417 the claims made. These sections clearly state the contributions made by this paper and are further
418 developed in the following sections.

419 Guidelines:

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

429 **2. Limitations**

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

431 Answer: **[Yes]**

432 Justification: This paper dives heavily into the limitations of our proposed optimizer, highlighting
433 the computational issues that come with it along with domains where it underperforms compared to
434 other optimizers. Additionally, we also discuss potential limitations that we did not come across in our
435 studies, but are plausible.

436 Guidelines:

- 437 • The answer NA means that the paper has no limitation while the answer No means that the paper
 438 has limitations, but those are not discussed in the paper.
 439 • The authors are encouraged to create a separate "Limitations" section in their paper.
 440 • The paper should point out any strong assumptions and how robust the results are to violations of
 441 these assumptions (e.g., independence assumptions, noiseless settings, model well-specification,
 442 asymptotic approximations only holding locally). The authors should reflect on how these
 443 assumptions might be violated in practice and what the implications would be.
 444 • The authors should reflect on the scope of the claims made, e.g., if the approach was only tested
 445 on a few datasets or with a few runs. In general, empirical results often depend on implicit
 446 assumptions, which should be articulated.
 447 • The authors should reflect on the factors that influence the performance of the approach. For
 448 example, a facial recognition algorithm may perform poorly when image resolution is low or
 449 images are taken in low lighting. Or a speech-to-text system might not be used reliably to provide
 450 closed captions for online lectures because it fails to handle technical jargon.
 451 • The authors should discuss the computational efficiency of the proposed algorithms and how
 452 they scale with dataset size.
 453 • If applicable, the authors should discuss possible limitations of their approach to address problems
 454 of privacy and fairness.
 455 • While the authors might fear that complete honesty about limitations might be used by reviewers
 456 as grounds for rejection, a worse outcome might be that reviewers discover limitations that
 457 aren't acknowledged in the paper. The authors should use their best judgment and recognize
 458 that individual actions in favor of transparency play an important role in developing norms that
 459 preserve the integrity of the community. Reviewers will be specifically instructed to not penalize
 460 honesty concerning limitations.

461 **3. Theory assumptions and proofs**

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

464 Answer: [Yes]

465 Justification: For each of our theoretical results, we give an extensive setup covering assumptions,
 466 lemmas, theorems, definitions, etc. Additionally, we validate each of these where necessary and
 467 provide sketches for the proofs as well.

468 Guidelines:

- 469 • The answer NA means that the paper does not include theoretical results.
 470 • All the theorems, formulas, and proofs in the paper should be numbered and cross-referenced.
 471 • All assumptions should be clearly stated or referenced in the statement of any theorems.
 472 • The proofs can either appear in the main paper or the supplemental material, but if they appear in
 473 the supplemental material, the authors are encouraged to provide a short proof sketch to provide
 474 intuition.
 475 • Inversely, any informal proof provided in the core of the paper should be complemented by
 476 formal proofs provided in appendix or supplemental material.
 477 • Theorems and Lemmas that the proof relies upon should be properly referenced.

478 **4. Experimental result reproducibility**

479 Question: Does the paper fully disclose all the information needed to reproduce the main experimental
 480 results of the paper to the extent that it affects the main claims and/or conclusions of the paper
 481 (regardless of whether the code and data are provided or not)?

482 Answer: [Yes]

483 Justification: This paper provides in-depth details and self-contained pseudocode for the algorithm we
 484 propose, along with hyperparameters, data splits, hardware details, a fixed random seed (42), and more
 485 for reproducibility. Additionally, we intend to release the source code used for the experimentation
 486 done in this study.

487 Guidelines:

- 488 • The answer NA means that the paper does not include experiments.
 489 • If the paper includes experiments, a No answer to this question will not be perceived well by the
 490 reviewers: Making the paper reproducible is important, regardless of whether the code and data
 491 are provided or not.
 492 • If the contribution is a dataset and/or model, the authors should describe the steps taken to make
 493 their results reproducible or verifiable.

- 494 • Depending on the contribution, reproducibility can be accomplished in various ways. For
 495 example, if the contribution is a novel architecture, describing the architecture fully might suffice,
 496 or if the contribution is a specific model and empirical evaluation, it may be necessary to either
 497 make it possible for others to replicate the model with the same dataset, or provide access to
 498 the model. In general, releasing code and data is often one good way to accomplish this, but
 499 reproducibility can also be provided via detailed instructions for how to replicate the results,
 500 access to a hosted model (e.g., in the case of a large language model), releasing of a model
 501 checkpoint, or other means that are appropriate to the research performed.
- 502 • While NeurIPS does not require releasing code, the conference does require all submissions
 503 to provide some reasonable avenue for reproducibility, which may depend on the nature of the
 504 contribution. For example
- 505 (a) If the contribution is primarily a new algorithm, the paper should make it clear how to
 506 reproduce that algorithm.
- 507 (b) If the contribution is primarily a new model architecture, the paper should describe the
 508 architecture clearly and fully.
- 509 (c) If the contribution is a new model (e.g., a large language model), then there should either be
 510 a way to access this model for reproducing the results or a way to reproduce the model (e.g.,
 511 with an open-source dataset or instructions for how to construct the dataset).
- 512 (d) We recognize that reproducibility may be tricky in some cases, in which case authors are
 513 welcome to describe the particular way they provide for reproducibility. In the case of
 514 closed-source models, it may be that access to the model is limited in some way (e.g.,
 515 to registered users), but it should be possible for other researchers to have some path to
 516 reproducing or verifying the results.

517 5. Open access to data and code

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

520 Answer: [Yes]

521 Justification: We submit our code as supplemental material, and intend to include a GitHub link in the
 522 final paper version if accepted for publication. This will include a README and requirements.txt
 523 so users can fully recreate our testing environment. This will provide commands for execution, data
 524 splits, etc.

525 Guidelines:

- 526 • The answer NA means that paper does not include experiments requiring code.
- 527 • Please see the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- 528 • While we encourage the release of code and data, we understand that this might not be possible,
 529 so “No” is an acceptable answer. Papers cannot be rejected simply for not including code, unless
 530 this is central to the contribution (e.g., for a new open-source benchmark).
- 531 • The instructions should contain the exact command and environment needed to run to reproduce
 532 the results. See the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- 533 • The authors should provide instructions on data access and preparation, including how to access
 534 the raw data, preprocessed data, intermediate data, and generated data, etc.
- 535 • The authors should provide scripts to reproduce all experimental results for the new proposed
 536 method and baselines. If only a subset of experiments are reproducible, they should state which
 537 ones are omitted from the script and why.
- 538 • At submission time, to preserve anonymity, the authors should release anonymized versions (if
 539 applicable).
- 540 • Providing as much information as possible in supplemental material (appended to the paper) is
 541 recommended, but including URLs to data and code is permitted.

544 6. Experimental setting/details

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

547 Answer: [Yes]

548 Justification: The main text and appendices give full details on data splits, hyperparameter selection,
 549 optimizer configurations, and training procedures. Additionally, the configuration and training files
 550 are submitted with this work.

551 Guidelines:

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

556 **7. Experiment statistical significance**

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

559 Answer: [Yes]

560 Justification: We report error bars in all plots, representing the standard error of the mean (SEM)
561 across independent runs, capturing variability due to random initialization and data splits. Error bars
562 and p-values (from Welch's t-test) are calculated as described in the appendix, with all methods and
563 assumptions clearly stated and referenced in the main text and figures.

564 Guidelines:

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

584 **8. Experiments compute resources**

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

587 Answer: [Yes]

588 Justification: For each experiment run, we record and report the type of compute resources used
589 (CPU/GPU model, memory, and execution time), as well as detailed system information and per-run
590 resource usage, which are included in the main text, appendices, and made available in the repository.
591 We also provide the compute requirements for individual runs and the total experiment, enabling full
592 reproducibility and transparency regarding computational resources.

593 Guidelines:

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

602 **9. Code of ethics**

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

605 Answer: [Yes]

606 Justification: The submission is fully anonymized and no ethical concerns arise from a foundational
607 optimizer study.

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

614 **10. Broader impacts**

615 Question: Does the paper discuss both potential positive societal impacts and negative societal impacts

616 of the work performed?

617 Answer: [NA]

618 Justification: This paper has no direct societal impact as it is foundational research, having the same

619 level of impact as any other optimizer.

620 Guidelines:

- 621 • The answer NA means that there is no societal impact of the work performed.
- 622 • If the authors answer NA or No, they should explain why their work has no societal impact or
- 623 why the paper does not address societal impact.
- 624 • Examples of negative societal impacts include potential malicious or unintended uses (e.g.,
- 625 disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deploy-
- 626 ment of technologies that could make decisions that unfairly impact specific groups), privacy
- 627 considerations, and security considerations.
- 628 • The conference expects that many papers will be foundational research and not tied to particular
- 629 applications, let alone deployments. However, if there is a direct path to any negative applications,
- 630 the authors should point it out. For example, it is legitimate to point out that an improvement in
- 631 the quality of generative models could be used to generate deepfakes for disinformation. On the
- 632 other hand, it is not needed to point out that a generic algorithm for optimizing neural networks
- 633 could enable people to train models that generate Deepfakes faster.
- 634 • The authors should consider possible harms that could arise when the technology is being used
- 635 as intended and functioning correctly, harms that could arise when the technology is being used
- 636 as intended but gives incorrect results, and harms following from (intentional or unintentional)
- 637 misuse of the technology.
- 638 • If there are negative societal impacts, the authors could also discuss possible mitigation strategies
- 639 (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitor-
- 640 ing misuse, mechanisms to monitor how a system learns from feedback over time, improving the
- 641 efficiency and accessibility of ML).

642 **11. Safeguards**

643 Question: Does the paper describe safeguards that have been put in place for responsible release of

644 data or models that have a high risk for misuse (e.g., pretrained language models, image generators, or

645 scraped datasets)?

646 Answer: [NA]

647 Justification: This paper poses no such risks, as this optimizer is foundational and does not stand to be

648 misused. It introduces no high-risk data or models requiring controlled release.

649 Guidelines:

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

658 **12. Licenses for existing assets**

659 Question: Are the creators or original owners of assets (e.g., code, data, models), used in the paper,

660 properly credited and are the license and terms of use explicitly mentioned and properly respected?

661 Answer: [Yes]

662 Justification: All third-party methods and datasets are properly cited with original references and

663 URLs, crediting the creators.

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

678 **13. New assets**

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

680 alongside the assets?

681 Answer: [Yes]

682 Justification: This paper introduces a new software asset, the Milo optimizer. Additionally, it provides

683 extensive documentation alongside it.

684 Guidelines:

- 685 • The answer NA means that the paper does not release new assets.
- 686 • Researchers should communicate the details of the dataset/code/model as part of their sub-
- 687 missions via structured templates. This includes details about training, license, limitations,
- 688 etc.
- 689 • The paper should discuss whether and how consent was obtained from people whose asset is
- 690 used.
- 691 • At submission time, remember to anonymize your assets (if applicable). You can either create an
- 692 anonymized URL or include an anonymized zip file.

693 **14. Crowdsourcing and research with human subjects**

694 Question: For crowdsourcing experiments and research with human subjects, does the paper include

695 the full text of instructions given to participants and screenshots, if applicable, as well as details about

696 compensation (if any)?

697 Answer: [NA]

698 Justification: This paper does not involve crowdsourcing nor research with human subjects. The nature

699 of our work does not involve it.

700 Guidelines:

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

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

709 Question: Does the paper describe potential risks incurred by study participants, whether such

710 risks were disclosed to the subjects, and whether Institutional Review Board (IRB) approvals (or an

711 equivalent approval/review based on the requirements of your country or institution) were obtained?

712 Answer: [NA]

713 Justification: This paper does not involve crowdsourcing nor research with human subjects. The nature

714 of our work does not involve it.

715 Guidelines:

- 716 • The answer NA means that the paper does not involve crowdsourcing nor research with human
- 717 subjects.

- 718 • Depending on the country in which research is conducted, IRB approval (or equivalent) may be
719 required for any human subjects research. If you obtained IRB approval, you should clearly state
720 this in the paper.
721 • We recognize that the procedures for this may vary significantly between institutions and
722 locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the guidelines for
723 their institution.
724 • For initial submissions, do not include any information that would break anonymity (if applica-
725 ble), such as the institution conducting the review.

726 **16. Declaration of LLM usage**

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

731 Answer: [NA]

732 Justification: LLMs are not an important, original, or non-standard component of our research. They
733 have no impact on our methodology other than training one in an experiment.

734 Guidelines:

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

739 Appendices

740 Contents

741 1 Introduction	1
742 2 Related Work	2
743 3 Algorithm	3
744 4 Convergence Guarantees for Milo	6
745 5 Experiments	7
746 6 Limitations and Future Work	9
747 7 Conclusion	9
748 Appendices	20
749 A Nomenclature	21
750 B Convergence Proof for Convex Settings	22
751 C Convergence Proof for Non-Convex Settings	26
752 D Experimentation	29

Table 2: Nomenclature

Symbol	Description
Model Parameters and Tensors	
p_i	Parameter tensor.
$\text{layer}(p)$	Function mapping a parameter p to its layer identifier.
$\text{vec}(\cdot)$	Operation that flattens a tensor into a vector.
Gradients and Grouping	
g_i	Gradient of parameter p_i .
$\nabla f(p)$	Gradient of the objective function f with respect to p .
$g_t \in \mathbb{R}^n$	Full (raw) gradient at iteration t (vector of dimension n).
n	Total number of gradient components.
m	Number of disjoint groups into which g_t is partitioned.
j	Group index, $j \in \{1, 2, \dots, m\}$.
G_j	j th group (subset) of gradient components from g_t .
$ G_j $	Number of components in group G_j .
μ_j	Empirical mean of the components in group G_j .
σ_j^2	Empirical variance of the components in group G_j .
ϵ	Stability constant: small positive constant added in normalization ($\sigma_j + \epsilon$) to avoid division by zero.
\tilde{g}_t	Normalized gradient at iteration t , defined as $[\tilde{g}_t]_i = \frac{[g_t]_i - \mu_j}{\sigma_j + \epsilon}$ for $i \in G_j$.
ρ_j	Variance contraction factor for group G_j , defined as $\rho_j = \frac{\sigma_j^2}{(\sigma_j + \epsilon)^2}$.
ρ	Global constant such that $\rho_j \leq \rho < 1$ for all groups.
Layer-Wise Grouping (if applicable)	
\mathcal{G}_l	Set of parameters (or corresponding gradient indices) in layer l .
v_l	Flattened gradient vector for layer l .
$ v_l $	Number of elements in v_l .
$v_l^{(j)}$	j th subgroup of v_l (when partitioning a layer's gradients).
$\mu_l^{(j)}$	Mean of subgroup j in layer l .
$\sigma_l^{(j)}$	Standard deviation of subgroup j in layer l .
\tilde{v}_l	Concatenated normalized gradient vector for layer l .
Optimization Variables and Algorithm Parameters	
w_t	Current model parameters (iterate) at iteration t .
f or F	Objective function (convex: f ; non-convex: F).
η	Step size (learning rate).
T	Total number of iterations.
D	Diameter of the feasible set \mathcal{X} .
\mathcal{X}	Feasible set of model parameters (subset of \mathbb{R}^d).
G	Uniform upper bound on the gradient norm.
Theoretical Analysis Parameters	
δ_t	Alignment error vector defined by $\tilde{g}_t = \nabla f(w_t) + \delta_t$.
ε	Alignment error bound: a scalar such that $ \langle \nabla f(w_t), \delta_t \rangle \leq \varepsilon \ \nabla f(w_t)\ ^2$. specifically, it is bounded as $ \langle \nabla f(w_t), \delta_t \rangle \leq \varepsilon \ \nabla f(w_t)\ ^2$, with $\varepsilon = \frac{\epsilon}{\sigma_{\min}} C(\tau)$.
α	Constant in $(0, 1)$ controlling the ratio between ϵ and σ_j to ensure proper scaling.
$C(\tau)$	Constant that depends on the sub-Gaussian parameter τ of the gradient components.
$\Theta(\cdot)$	Tight asymptotic bound.
$O(\cdot)$	Big-O notation.
Regret Analysis	
$R(T)$	Cumulative regret over T iterations.
$R_{SGD}(T)$	Regret bound for standard SGD.
$R_{Milo}(T)$	Regret bound for Normalized SGD with Layer-Wise Grouping (Milo).

754 **B Convergence Proof for Convex Settings**

755 In this section we present the main convergence proof. We begin by introducing the definitions, assumptions,
 756 and lemmas that will be used throughout.

757 **B.1 Proof Sketch of Theorem**

758 Under standard L -smoothness, bounded-gradient, and variance-contraction/moment conditions, we obtain the
 759 following.

760 *Proof Sketch of Theorem B.* We outline the main steps under Assumptions B.1–B.4 and Lemmas B.4, B.2.

1. Smoothness Step. By L -smoothness,

$$f(w_{t+1}) \leq f(w_t) + \langle \nabla f(w_t), w_{t+1} - w_t \rangle + \frac{L}{2} \|w_{t+1} - w_t\|^2.$$

Using the update $w_{t+1} = w_t - \eta \tilde{g}_t$ gives

$$f(w_{t+1}) \leq f(w_t) - \eta \langle \nabla f(w_t), \tilde{g}_t \rangle + \frac{L\eta^2}{2} \|\tilde{g}_t\|^2.$$

2. Variance-Contraction Bound. By Assumption B.3 and Lemma B.4, each group satisfies $\|[\tilde{g}_t]_{G_j}\|^2 \leq \rho |G_j|$. Summing yields

$$\|\tilde{g}_t\|^2 \leq \rho n.$$

3. Alignment Control. Decompose $\tilde{g}_t = \nabla f(w_t) + \delta_t$. Lemma B.2 gives

$$|\langle \nabla f(w_t), \delta_t \rangle| \leq \varepsilon \|\nabla f(w_t)\|^2,$$

so

$$\langle \nabla f(w_t), \tilde{g}_t \rangle \geq (1 - \varepsilon) \|\nabla f(w_t)\|^2.$$

4. One-Step Progress. Combine the above bounds:

$$f(w_{t+1}) \leq f(w_t) - \eta(1 - \varepsilon) \|\nabla f(w_t)\|^2 + \frac{L\eta^2}{2} \rho n.$$

5. Telescoping Sum. Summing over $t = 1 \dots T$ and using $f(w_{T+1}) \geq f(w^*)$:

$$\sum_{t=1}^T \|\nabla f(w_t)\|^2 \leq \frac{f(w_1) - f(w^*)}{\eta(1 - \varepsilon)} + \frac{L\eta T}{2(1 - \varepsilon)} \rho n.$$

6. Learning-Rate Choice. Setting $\eta = \Theta(1/\sqrt{T})$ balances the two terms and yields

$$\frac{1}{T} \sum_{t=1}^T \|\nabla f(w_t)\|^2 = O(T^{-1/2}).$$

7. Regret Comparison. Standard SGD regret is $O(GD\sqrt{T})$. Replacing G by $\sqrt{\rho n}$ gives Milo

$$R_{\text{Milo}}(T) = O(D\sqrt{\rho n T}) < O(DG\sqrt{T}) = R_{\text{SGD}}(T),$$

761 showing improved dependence on the variance.

762

763 **B.2 Definitions and Assumptions**

764 **Definition B.1** (Convex Function). A function $f : \mathbb{R}^d \rightarrow \mathbb{R}$ is convex if for all $x, y \in \mathbb{R}^d$ and all $\lambda \in [0, 1]$,

$$\lambda f(x) + (1 - \lambda)f(y) \geq f(\lambda x + (1 - \lambda)y). \quad (\text{B.1})$$

765 **Remark B.1.** Notice that for a convex function f , one can lower bound f by a hyperplane at any point of its
 766 domain. In particular, if f is differentiable, then for all $x, y \in \mathbb{R}^d$,

$$f(y) \geq f(x) + \nabla f(x)^T (y - x). \quad (\text{B.2})$$

767 **Assumption B.1** (Smoothness). *The function $f : \mathbb{R}^d \rightarrow \mathbb{R}$ is differentiable and L -smooth; that is, for all
768 $x, y \in \mathbb{R}^d$,*

$$\|\nabla f(x) - \nabla f(y)\|_2 \leq L\|x - y\|_2. \quad (\text{B.3})$$

769 *Often, f is given as a finite sum $f(x) = \frac{1}{n} \sum_{i=1}^n f_i(x)$ with individual smoothness constants L_i ; we denote
770 $L_{\max} = \max\{L_1, \dots, L_n\}$.*

771 [See, e.g., [Nesterov, 1983] for the smoothness condition.]

772 **Assumption B.2** (Bounded Gradient Norm). *For every $x \in \mathcal{X}$ (the feasible set), the gradient is bounded:*

$$\|\nabla f(x)\|_2 \leq G, \quad (\text{B.4})$$

773 *with $G > 0$.*

774 [This is standard in stochastic optimization; see, e.g., [Bottou et al., 2018].]

775 **Assumption B.3** (Variance Contraction via Normalization). *Let $g_t = \nabla f(w_t)$ be the gradient at iteration t .
776 Partition g_t into disjoint groups $\{G_j\}_{j=1}^m$ (with $\sum_{j=1}^m |G_j| = n$). For each group G_j , define the empirical
777 mean and variance as*

$$\mu_j = \frac{1}{|G_j|} \sum_{i \in G_j} [g_t]_i, \quad \sigma_j^2 = \frac{1}{|G_j|} \sum_{i \in G_j} ([g_t]_i - \mu_j)^2. \quad (\text{B.5})$$

778 *The normalized gradient is then defined by*

$$[\tilde{g}_t]_i = \frac{[g_t]_i - \mu_j}{\sigma_j + \epsilon}, \quad i \in G_j, \quad (\text{B.6})$$

779 *Because the mean is subtracted in Equation (B.6) (centering the data), the variance of the normalized gradient
780 in G_j becomes*

$$\tilde{\sigma}_j^2 = \frac{\sigma_j^2}{(\sigma_j + \epsilon)^2}. \quad (\text{B.7})$$

781 *Define the variance contraction factor as*

$$\rho_j = \frac{\sigma_j^2}{(\sigma_j + \epsilon)^2}. \quad (\text{B.8})$$

782 *We assume that there exists a global constant $\rho \in [0, 1)$ such that*

$$\rho_j \leq \rho, \quad \forall j. \quad (\text{B.9})$$

783 [See, e.g., [Wu and He, 2018] for the effects of group normalization on variance contraction.]

784 **Assumption B.4** (Moment Conditions). *For each parameter tensor p with flattened gradient $v = \text{vec}(g(p))$,
785 assume:*

786 1. **(Large Group Size)** *Each group G_j (obtained by partitioning v) satisfies $|G_j| \geq N_0$ so that the
787 empirical estimates*

$$\mu_j = \frac{1}{|G_j|} \sum_{i \in G_j} v_i, \quad \sigma_j^2 = \frac{1}{|G_j|} \sum_{i \in G_j} (v_i - \mu_j)^2 \quad (\text{B.10})$$

788 *concentrate around the true moments.*

789 2. **(Bounded Higher-Order Moments)** *There exists a constant M_k and some $k \geq 3$ such that*

$$\mathbb{E}[|v_i - \mu_j|^k] \leq M_k. \quad (\text{B.11})$$

790 3. **(Symmetry)** *The distribution of v_i for $i \in G_j$ is symmetric (or nearly so) around μ_j .*

791 [Moment conditions as in [Vershynin, 2018].]

792 **B.3 Auxiliary Lemmas**

793 **Lemma B.1** (Concentration of Empirical Means). *Let v_1, \dots, v_N be independent random variables in a group
794 G_j with true mean μ and assume $|v_i - \mu| \leq B$ almost surely. Then, for any $\delta > 0$,*

$$\Pr\left(\left|\frac{1}{N} \sum_{i=1}^N v_i - \mu\right| \geq \delta\right) \leq 2 \exp\left(-\frac{2N\delta^2}{B^2}\right). \quad (\text{B.12})$$

795 A similar bound holds for the empirical variance.

796 [This is an application of Hoeffding's inequality; see [Hoeffding, 1963].]

797 **Lemma B.2** (Formal Alignment of Normalized Gradient). *Assume that for a given group G_j the gradient
798 components can be expressed as*

$$v_i = \mu_j + \sigma_j z_i, \quad i \in G_j, \quad (\text{B.13})$$

799 where z_i are independent sub-Gaussian random variables with zero mean, unit variance, and sub-Gaussian
800 parameter τ . Define the normalized elements by

$$\tilde{v}_i = \frac{v_i - \mu_j}{\sigma_j + \epsilon}. \quad (\text{B.14})$$

801 Then,

$$\tilde{v}_i = \frac{\sigma_j}{\sigma_j + \epsilon} z_i. \quad (\text{B.15})$$

802 Define the scaling factor $P_j := \frac{\sigma_j}{\sigma_j + \epsilon}$. If for a fixed $\alpha \in (0, 1)$ it holds that $\epsilon \leq \alpha \sigma_j$ (i.e. $P_j \geq \frac{1}{1+\alpha}$), then

$$1 - P_j \leq \frac{\epsilon}{\sigma_j}. \quad (\text{B.16})$$

803 Now, writing the decomposition

$$\tilde{g}_t = \nabla f(w_t) + \delta_t, \quad (\text{B.17})$$

804 the above results (with appropriate concentration bounds) imply that

$$|\langle \nabla f(w_t), \delta_t \rangle| \leq \epsilon \|\nabla f(w_t)\|^2, \quad (\text{B.18})$$

805 with $\epsilon = \frac{\epsilon}{\sigma_{\min}} C(\tau)$ (here $\sigma_{\min} = \min_j \sigma_j$ and $C(\tau)$ is a constant depending on τ).

806 [See, e.g., [Duchi et al., 2011] for related analysis on adaptive and normalized gradient methods.]

807 **Lemma B.3** (Robustness for Small σ_j with Explicit ϵ). *Assume that for each group G_j either:*

808 1. $\sigma_j \geq \sigma_{\min} > 0$, or

809 2. If $\sigma_j < \sigma_{\min}$, then $|v_i - \mu_j| \leq \beta \sigma_{\min}$ for all $i \in G_j$,

810 for some constants $\sigma_{\min} > 0$ and $\beta > 0$. Then, by choosing

$$\epsilon \leq \gamma \sigma_{\min}, \quad (\text{B.19})$$

811 for a given $\gamma \in (0, 1)$, it follows that for every $i \in G_j$,

$$|[\tilde{g}_t]_i| \leq \frac{\beta}{\gamma}. \quad (\text{B.20})$$

812 **Lemma B.4** (Bounded Variance Contraction). *Let $\epsilon > 0$ be such that for every group G_j (and for all iterations)
813 it holds that*

$$\epsilon \geq c \sigma_j, \quad (\text{B.21})$$

814 for some constant $c > 0$. Then, the variance contraction factor satisfies

$$\rho_j = \frac{\sigma_j^2}{(\sigma_j + \epsilon)^2} \leq \left(\frac{1}{1+c}\right)^2 < 1. \quad (\text{B.22})$$

815 [This result is motivated by normalization techniques in [Ioffe and Szegedy, 2015].]

816 **B.4 Proof of Convergence for Convex Settings**

817 We now proceed with the convergence proof. The update rule for Milo is given (after normalization) by

$$w_{t+1} = w_t - \eta \tilde{g}_t, \quad (\text{B.23})$$

818 where η is the learning rate and the normalized gradient \tilde{g}_t is defined in eq. Equation (B.6). Throughout the
819 proof we also use the following decomposition:

$$\tilde{g}_t = \nabla f(w_t) + \delta_t, \quad (\text{B.24})$$

820 where δ_t represents the alignment error (as introduced in Lemma B.2).

821 **Step 1: Smoothness Inequality.** By Assumption B.1 (Smoothness), for any $x, y \in \mathbb{R}^d$ we have

$$f(y) \leq f(x) + \langle \nabla f(x), y - x \rangle + \frac{L}{2} \|y - x\|^2. \quad (\text{B.25})$$

822 Taking $x = w_t$ and $y = w_{t+1} = w_t - \eta \tilde{g}_t$ yields

$$f(w_{t+1}) \leq f(w_t) - \eta \langle \nabla f(w_t), \tilde{g}_t \rangle + \frac{L\eta^2}{2} \|\tilde{g}_t\|^2. \quad (\text{B.26})$$

823 Here:

- 824 • $\nabla f(w_t)$ is the true gradient at w_t ,
- 825 • \tilde{g}_t is the normalized gradient from Equation (B.6), and
- 826 • $\langle \nabla f(w_t), \tilde{g}_t \rangle$: The inner product measuring the alignment between the true gradient and the normalized gradient.
- 828 • L : The Lipschitz constant of ∇f , governing the curvature of f , as in Assumption B.1.

829 **Step 2: Bounding the Norm of \tilde{g}_t .** Within each group G_j , by the variance contraction property in
830 Assumption B.3 and Lemma B.4,

$$\|[\tilde{g}_t]_{G_j}\|^2 = \frac{|G_j| \sigma_j^2}{(\sigma_j + \epsilon)^2} \leq \rho |G_j|. \quad (\text{B.27})$$

831 Summing over all groups, we obtain

$$\|\tilde{g}_t\|^2 \leq \sum_j \rho |G_j| = \rho n. \quad (\text{B.28})$$

832 **Step 3: Handling the Alignment Error.** We decompose the normalized gradient as

$$\tilde{g}_t = \nabla f(w_t) + \delta_t. \quad (\text{B.29})$$

833 Then, by Lemma B.2 (Formal Alignment of Normalized Gradient), under the moment conditions (Assumption
834 B.4) and the condition $\epsilon \leq \alpha \sigma_j$ for all j , we have

$$|\langle \nabla f(w_t), \delta_t \rangle| \leq \varepsilon \|\nabla f(w_t)\|^2, \quad (\text{B.30})$$

835 with $\varepsilon = \frac{\epsilon}{\sigma_{\min}} C(\tau)$. Consequently,

$$\langle \nabla f(w_t), \tilde{g}_t \rangle \geq (1 - \varepsilon) \|\nabla f(w_t)\|^2. \quad (\text{B.31})$$

836 **Step 4: Inserting the Bounds.** Substitute the alignment bound and the norm bound into the smoothness
837 inequality:

$$f(w_{t+1}) \leq f(w_t) - \eta(1 - \varepsilon) \|\nabla f(w_t)\|^2 + \frac{L\eta^2}{2} \rho n. \quad (\text{B.32})$$

838 **Step 5: Telescoping.** Summing from $t = 1$ to T yields

$$\sum_{t=1}^T [f(w_t) - f(w_{t+1})] \geq \eta(1 - \varepsilon) \sum_{t=1}^T \|\nabla f(w_t)\|^2 - \frac{L\eta^2 T}{2} \rho n. \quad (\text{B.33})$$

839 Since the sum telescopes to $f(w_1) - f(w_{T+1})$ and $f(w_{T+1}) \geq f(w^*)$ (where w^* is an optimal point), we have

$$f(w_1) - f(w^*) \geq \eta(1 - \varepsilon) \sum_{t=1}^T \|\nabla f(w_t)\|^2 - \frac{L\eta^2 T}{2} \rho n. \quad (\text{B.34})$$

840 Dividing both sides by $\eta(1 - \varepsilon)T$ gives

$$\frac{1}{T} \sum_{t=1}^T \|\nabla f(w_t)\|^2 \leq \frac{f(w_1) - f(w^*)}{\eta(1 - \varepsilon)T} + \frac{L\eta}{2(1 - \varepsilon)} \rho n. \quad (\text{B.35})$$

841 **Step 6: Choosing the Learning Rate.** By setting

$$\eta = \Theta\left(\frac{1}{\sqrt{T}}\right), \quad (\text{B.36})$$

842 the two terms on the right-hand side are balanced, leading to an overall convergence rate of

$$\frac{1}{T} \sum_{t=1}^T \|\nabla f(w_t)\|^2 = O\left(\frac{1}{\sqrt{T}}\right). \quad (\text{B.37})$$

843 **Step 7: Regret Bound.** Standard regret analysis for SGD yields

$$R_{\text{SGD}}(T) \leq \frac{D^2}{2\eta} + \frac{\eta G^2 T}{2}. \quad (\text{B.38})$$

844 For Milo, if we denote by \tilde{G} the effective bound on $\|\tilde{g}_t\|$ (with $\tilde{G} \leq \sqrt{\rho n}$), then the regret can be bounded as

$$R_{\text{Milo}}(T) \leq \frac{D^2}{2\eta} + \frac{\eta \tilde{G}^2 T}{2}. \quad (\text{B.39})$$

845 Choosing $\eta = \frac{D}{\tilde{G}\sqrt{T}}$ gives

$$R_{\text{Milo}}(T) \leq D \tilde{G} \sqrt{T} < R_{\text{SGD}}(T) \leq D G \sqrt{T}, \quad (\text{B.40})$$

846 thus demonstrating the advantage of variance reduction.

847 B.5 Discussion and Practical Relevance

848 The final convergence bound explicitly depends on the alignment error ε , the variance contraction factor ρ , and
849 the learning rate η :

$$\frac{1}{T} \sum_{t=1}^T \|\nabla f(w_t)\|^2 \leq \frac{f(w_1) - f(w^*)}{\eta(1-\varepsilon)T} + \frac{L\eta}{2(1-\varepsilon)} \rho n. \quad (\text{B.41})$$

850 By balancing the two terms with $\eta = \Theta(1/\sqrt{T})$, this yields the stated $O(1/\sqrt{T})$ convergence rate. Moreover,
851 replacing the usual gradient bound G by $\sqrt{\rho n}$ in the standard SGD regret analysis gives

$$R_{\text{Milo}}(T) = O(D\sqrt{\rho n T}) < O(DG\sqrt{T}) = R_{\text{SGD}}(T),$$

852 demonstrating Milo's tighter dependence on variance.

853 Lemmas B.1 and B.2 show that under the moment conditions (Assumption B.4) and for sufficiently large group
854 sizes, the empirical estimates μ_j, σ_j^2 concentrate around their true values and the normalized gradient remains
855 well aligned with $\nabla f(w_t)$, yielding an alignment error $\varepsilon = O(\frac{\epsilon}{\sigma_{\min}})$. Lemma B.4 guarantees a global variance
856 contraction $\rho < 1$ whenever $\epsilon \geq c\sigma_j$, and Lemma B.3 ensures that even in low-variance regimes ($\sigma_j < \sigma_{\min}$)
857 the updates stay bounded by choosing ϵ relative to σ_{\min} .

858 By making these parameter dependencies explicit, our analysis directly informs practical implementation:

- 859 1. Choose group sizes large enough to invoke concentration (Assumption B.4).
- 860 2. Select $\epsilon = O(\sigma_{\min})$ to control both alignment error ε and variance contraction ρ .
- 861 3. Verify sub-Gaussian or bounded-moment behavior of gradient components to satisfy the lemmas.

862 Empirical results (cf. Section 5) confirm that in realistic deep-learning settings these constants remain small,
863 yielding smoother loss trajectories, reduced sensitivity to learning-rate tuning, and better generalization. This
864 completes our self-contained convex convergence proof and highlights Milo's theoretical and practical advantages
865 over standard SGD.

866 B.5.1 Practical Remarks.

867 In practice, unbiasedness holds up to $O(\epsilon/\sigma_{\min})$ when group sizes are large and gradient-component distributions
868 are near-symmetric; variance is contracted by $\rho_j \leq (\sigma_j/(\sigma_j + \epsilon))^2 < 1$ whenever $\epsilon \geq c\sigma_j$; and even if some
869 σ_j are tiny, choosing $\epsilon = O(\sigma_{\min})$ (Lemma B.3) keeps updates bounded.

870 C Convergence Proof for Non-Convex Settings

871 In this section, we present a detailed convergence analysis for Milo when the objective function is *non-convex*.
872 Although nonconvexity precludes guarantees of reaching a global optimum, under standard smoothness and
873 boundedness assumptions, we can establish a sublinear rate in terms of the (squared) gradient norm.

874 *Proof Sketch of Theorem C.* Under Assumptions C.1, C.2, and the variance contraction/moment conditions
875 (Assumptions B.3, B.4), we outline the main steps.

1. Descent Lemma. By L -smoothness of F , for $w_{t+1} = w_t - \eta \tilde{g}_t$:

$$F(w_{t+1}) \leq F(w_t) - \eta \nabla F(w_t), \tilde{g}_t + \frac{L\eta^2}{2} \|\tilde{g}_t\|^2.$$

2. Alignment Bound. Decompose $\tilde{g}_t = \nabla F(w_t) + \delta_t$. From Lemma B.2,

$$\nabla F(w_t), \tilde{g}_t \geq (1 - \varepsilon) \|\nabla F(w_t)\|^2.$$

3. Norm Control. Using variance contraction (Lemma B.4),

$$\|\tilde{g}_t\|^2 \leq \rho n.$$

4. One-Step Inequality. Combine the above into the descent:

$$F(w_{t+1}) \leq F(w_t) - \eta(1 - \varepsilon) \|\nabla F(w_t)\|^2 + \frac{L\eta^2}{2} \rho n.$$

Rearrange:

$$\eta(1 - \varepsilon) \|\nabla F(w_t)\|^2 \leq F(w_t) - F(w_{t+1}) + \frac{L\eta^2}{2} \rho n.$$

5. Telescoping Sum. Summing for $t = 0 \dots T - 1$ and using $F(w_T) \geq F^*$:

$$\sum_{t=0}^{T-1} \|\nabla F(w_t)\|^2 \leq \frac{F(w_0) - F^*}{\eta(1 - \varepsilon)} + \frac{L\eta T}{2(1 - \varepsilon)} \rho n.$$

6. Rate via Step-Size. Choose $\eta = \min\{1/L, c/\sqrt{T}\}$ for constant $c > 0$. Then both terms scale as $O(T^{-1/2})$, yielding

$$\min_{t < T} \|\nabla F(w_t)\|^2 = O(T^{-1/2}).$$

876

□

877 C.1 Assumptions and Preliminaries

878 Let $F : \mathbb{R}^d \rightarrow \mathbb{R}$ be a (possibly) non-convex function. In addition to the variance reduction ingredients from the
879 convex analysis, we assume the following:

880 **Assumption C.1** (Smoothness for Non-Convex Functions). *F is differentiable and L-smooth; that is, for all
881 $x, y \in \mathbb{R}^d$,*

$$\|\nabla F(x) - \nabla F(y)\|_2 \leq L\|x - y\|_2. \quad (\text{C.1})$$

882 **Assumption C.2** (Bounded Gradient). *For every $x \in \mathbb{R}^d$, the gradient is bounded:*

$$\|\nabla F(x)\|_2 \leq G, \quad (\text{C.2})$$

883 with $G > 0$.

884 We also assume that the grouping and normalization procedure is defined exactly as in the convex case. That is,
885 for iteration t we have

$$g_t = \nabla F(w_t), \quad (\text{C.3})$$

886 which is partitioned into groups $\{G_j\}_{j=1}^m$ with

$$\mu_j = \frac{1}{|G_j|} \sum_{i \in G_j} [g_t]_i, \quad \sigma_j^2 = \frac{1}{|G_j|} \sum_{i \in G_j} ([g_t]_i - \mu_j)^2, \quad (\text{C.4})$$

887 and the normalized gradient is defined as

$$[\tilde{g}_t]_i = \frac{[g_t]_i - \mu_j}{\sigma_j + \epsilon}, \quad i \in G_j. \quad (\text{C.5})$$

888 We assume that Assumptions B.3 and B.4 from the convex proof (regarding variance contraction via normalization
889 and moment conditions) hold unchanged here.

890 C.2 Proof of Convergence for Non-Convex Settings

891 The Milo update is given by

$$w_{t+1} = w_t - \eta \tilde{g}_t. \quad (\text{C.6})$$

892 Since F is L -smooth (Assumption C.1), we have the descent lemma:

$$F(w_{t+1}) \leq F(w_t) + \langle \nabla F(w_t), w_{t+1} - w_t \rangle + \frac{L}{2} \|w_{t+1} - w_t\|_2^2. \quad (\text{C.7})$$

893 Substituting $w_{t+1} = w_t - \eta \tilde{g}_t$ into Equation (C.7), we obtain

$$F(w_{t+1}) \leq F(w_t) - \eta \langle \nabla F(w_t), \tilde{g}_t \rangle + \frac{L\eta^2}{2} \|\tilde{g}_t\|_2^2. \quad (\text{C.8})$$

894 **Decomposition and Alignment.** We write the normalized gradient as

$$\tilde{g}_t = \nabla F(w_t) + \delta_t, \quad (\text{C.9})$$

895 where δ_t is the error due to the nonlinearity of the normalization. Under the moment conditions (Assumption B.4)
896 and by using Lemma B.2 (Formal Alignment of Normalized Gradient), we can bound the deviation:

$$\langle \nabla F(w_t), \tilde{g}_t \rangle \geq (1 - \varepsilon) \|\nabla F(w_t)\|_2^2, \quad (\text{C.10})$$

897 for some ε (dependent on ϵ/σ_{\min}).

898 **Bounding the Norm of \tilde{g}_t .** As in the convex setting (see Lemma B.4), the group normalization yields

$$\|\tilde{g}_t\|_2^2 \leq \rho n, \quad (\text{C.11})$$

899 where n is the total number of parameters and $\rho \in [0, 1)$ is the variance contraction factor.

900 **Combining the Bounds.** Substitute the alignment bound Equation (C.10) and the norm bound Equation (C.11) into Equation (C.8):

$$F(w_{t+1}) \leq F(w_t) - \eta(1 - \varepsilon) \|\nabla F(w_t)\|_2^2 + \frac{L\eta^2}{2} \rho n. \quad (\text{C.12})$$

902 Rearrange Equation (C.12) to obtain an inequality on the gradient norm:

$$\eta(1 - \varepsilon) \|\nabla F(w_t)\|_2^2 \leq F(w_t) - F(w_{t+1}) + \frac{L\eta^2}{2} \rho n. \quad (\text{C.13})$$

903 **Telescoping over Iterations.** Summing Equation (C.13) over $t = 0, \dots, T - 1$ results in:

$$\eta(1 - \varepsilon) \sum_{t=0}^{T-1} \|\nabla F(w_t)\|_2^2 \leq F(w_0) - F(w_T) + \frac{L\eta^2}{2} \rho nT. \quad (\text{C.14})$$

904 Since $F(w_T) \geq F^*$ (with $F^* = \min_w F(w)$ or a known lower bound), we have:

$$\eta(1 - \varepsilon) \sum_{t=0}^{T-1} \|\nabla F(w_t)\|_2^2 \leq F(w_0) - F^* + \frac{L\eta^2}{2} \rho nT. \quad (\text{C.15})$$

905 Dividing both sides of Equation (C.15) by $\eta(1 - \varepsilon)T$ yields:

$$\frac{1}{T} \sum_{t=0}^{T-1} \|\nabla F(w_t)\|_2^2 \leq \frac{F(w_0) - F^*}{\eta(1 - \varepsilon)T} + \frac{L\eta\rho n}{2(1 - \varepsilon)}. \quad (\text{C.16})$$

906 **Choosing the Learning Rate.** By selecting the step-size as

$$\eta = \min\left\{\frac{1}{L}, \frac{c}{\sqrt{T}}\right\}, \quad (\text{C.17})$$

907 for some constant $c > 0$, the two terms on the right-hand side of Equation (C.16) balance and both decay as
908 $O(1/\sqrt{T})$. Consequently, we obtain

$$\min_{t=0, \dots, T-1} \|\nabla F(w_t)\|_2^2 = O\left(\frac{1}{\sqrt{T}}\right), \quad (\text{C.18})$$

909 which establishes that Milo converges (in terms of the squared gradient norm) at a sublinear rate despite
910 nonconvexity.

911 C.3 Summary of the Proof

912 1. **Smoothness:** Using L -smoothness we derived

$$F(w_{t+1}) \leq F(w_t) - \eta \langle \nabla F(w_t), \tilde{g}_t \rangle + \frac{L\eta^2}{2} \|\tilde{g}_t\|_2^2.$$

913 2. **Alignment:** The error decomposition $\tilde{g}_t = \nabla F(w_t) + \delta_t$ combined with Lemma B.2 gives

$$\langle \nabla F(w_t), \tilde{g}_t \rangle \geq (1 - \varepsilon) \|\nabla F(w_t)\|_2^2.$$

914 3. **Norm Bound:** Lemma B.4 ensures that $\|\tilde{g}_t\|_2^2 \leq \rho n$.

915 4. **Telescoping:** Summing over iterations and rearranging leads to

$$\frac{1}{T} \sum_{t=0}^{T-1} \|\nabla F(w_t)\|_2^2 \leq \frac{F(w_0) - F^*}{\eta(1 - \varepsilon)T} + \frac{L\eta\rho n}{2(1 - \varepsilon)}.$$

916 5. **Rate:** With the choice of step-size in Equation (C.17), we conclude that

$$\min_{t=0, \dots, T-1} \|\nabla F(w_t)\|_2^2 = O\left(\frac{1}{\sqrt{T}}\right).$$

917 **C.4 Discussion and Practical Implications**

918 The preceding analysis shows that even when F is non-convex, Milo's normalized gradients drive a controlled
 919 descent of the objective and guarantee

$$\min_{t=0,\dots,T-1} \|\nabla F(w_t)\|_2^2 = O(T^{-1/2}).$$

920 The key technical ingredients are:

- 921 • **Smoothness:** L -smoothness ensures F does not change too abruptly, yielding the descent lemma.
- 922 • **Variance Contraction:** Group-wise normalization (Assumption B.3 and Lemma B.4) tightly controls
 923 $\|\tilde{g}_t\|$.
- 924 • **Alignment Control:** Lemma B.2 bounds the misalignment δ_t so that $\nabla F(w_t), \tilde{g}_t \geq (1 - \varepsilon) \|\nabla F(w_t)\|^2$.

926 By choosing $\eta = \min\{1/L, c/\sqrt{T}\}$, the telescoped one-step bounds give the claimed $O(T^{-1/2})$ rate in the
 927 squared gradient norm. Crucially, this leverages the same normalization strategy that yields variance contraction
 928 in the convex case, now coupled with careful alignment control in the non-convex regime.

929 This result provides rigorous theoretical backing for Milo's empirical success across both convex and non-convex
 930 problems, explaining why normalized updates can both stabilize training and accelerate convergence even on
 931 highly non-convex objectives.

932 **D Experimentation**

933 **D.1 Experiment Hardware**

Table 3: Static Hardware Configuration

Property	Value
Platform	Windows-11-10.0.22631-SP0
Processor	AMD Ryzen 9 9900X
Physical CPU cores	12
Logical CPU threads	24
Total system memory	61.6 GB
Disk capacity used	931.4 GB
CUDA available	True
Compute type	GPU
GPU count	1
GPU model	NVIDIA GeForce RTX 5070 Ti,
GPU memory	16.0 GB

934 **D.2 Libraries Used**

935 We rely on the following open-source software (version, license, URL, citation):

- 936 • **PyTorch** (v2.6.0; BSD 3-clause License; <https://pytorch.org>) [Paszke et al., 2019]
- 937 • **torchvision** (v0.21.0; BSD 3-clause License; <https://github.com/pytorch/vision>) [PyTorch
 938 Contributors, 2023]
- 939 • **seaborn** (v0.13.2; BSD License; <https://seaborn.pydata.org>) [Waskom, 2021]
- 940 • **matplotlib** (v3.10.0; PSF License; <https://matplotlib.org>) [Hunter, 2007]
- 941 • **NumPy** (v2.1.3; BSD License; <https://numpy.org>) [van der Walt et al., 2011]
- 942 • **NovoGrad** optimizer (preprint; Apache-2.0; <https://arxiv.org/abs/1905.11286>) [Ginsburg
 943 et al., 2019]
- 944 • **pandas** (v2.2.3; BSD License; <https://pandas.pydata.org>) [McKinney, 2010]
- 945 • **JSON module** (Python 3.10 stdlib; PSF License; <https://docs.python.org/3/library/json.html>) [Crockford, 2006]

- **Optuna** (v4.3.0; MIT License; <https://optuna.org>) [Akiba et al., 2019]
- **scikit-learn** (v1.6.1; BSD License; <https://scikit-learn.org>) [Pedregosa et al., 2011]
- **SciPy** (v1.15.1; BSD License; <https://scipy.org>) [Virtanen et al., 2020]
- **itertools** (Python 3.10 stdlib; PSF License; <https://docs.python.org/3/library/itertools.html>) [Python Software Foundation, 2023]
- **Gymnasium** (v1.1.0; MIT License; <https://gymnasium.farama.org>) [Brockman et al., 2016]
- **Transformers** (v4.30.0; Apache-2.0; <https://github.com/huggingface/transformers>) [Wolf et al., 2020]
- **train-lm-from-scratch** (MIT License; <https://github.com/FareedKhan-dev/train-lm-from-scratch>) [FareedKhan-dev, 2025]
- **NovoGrad-pytorch** (MIT License; <https://github.com/lonePatient/NovoGrad-pytorch>) [lonePatient, 2019]

959 D.3 Datasets

960 We conduct all experiments on the following publicly released datasets:

- **MNIST** (Public Domain; <http://yann.lecun.com/exdb/mnist/>) [LeCun et al., 1998]
- **CIFAR-10** (2009 Technical Report; MIT License; <https://www.cs.toronto.edu/~kriz/cifar.html>) [Krizhevsky and Hinton, 2009]
- **CIFAR-100** (2009 Technical Report; MIT License; <https://www.cs.toronto.edu/~kriz/cifar.html>) [Krizhevsky and Hinton, 2009]
- **The Pile** (Apache 2.0 License; <https://pile.eleuther.ai/>) [Gao et al., 2020]

967 D.4 Error Bar Computation

968 Error bars in our experiments quantify the uncertainty in each metric across multiple runs. They are computed as
969 follows:

970 1. Calculation

- **Standard Error of the Mean (SEM):**

$$\text{SEM} = \frac{\text{SD}}{\sqrt{n}}$$

972 where SD is the sample standard deviation and n is the number of independent runs.

- **95% Confidence Interval (CI):**

$$\text{CI} = t_{0.975, n-1} \times \text{SEM},$$

974 with $t_{0.975, n-1}$ the two-sided critical value from the Student's t -distribution.

- For each epoch or step, we aggregate the metric values over n runs, compute the mean, SD, SEM, and then plot the mean \pm CI as error bands.

977 2. Sources of Variability

- *Run-to-run stochasticity*: random initialization, data shuffling.
- *Optimizer/model differences*: varying update rules or architectures.
- *Data subsampling*: different train/validation/test splits.
- *Hardware noise*: floating-point precision and device variability.

982 3. Assumptions

- **Normality**: metric distributions are approximately Gaussian.
- **Independence**: each run is statistically independent.
- **Sufficient sample size**: $n > 1$ for reliable SEM and CI estimates.
- **Consistency**: identical hyperparameters and data splits across runs.

987 **4. Implementation** The function `calculate_statistics()` computes mean, SD, SEM, and CI for each
 988 metric, and `plot_seaborn_style_with_error_bars()` renders the resulting error bands in the figures.

989 D.5 Experiment 1: Loss Map Predictions

990 The objective of this experiment is to evaluate multiple optimizers on different loss maps to provide an
 991 understanding of their behavior in complex loss environments. In each optimization problem, the optimizers
 992 were run 5 times and their outputs were averaged, however, since it is a static environment there was never any
 993 variation in their behavior.

994 D.5.1 Himmelblau Function and Its Gradient

995 The Himmelblau function is defined as:

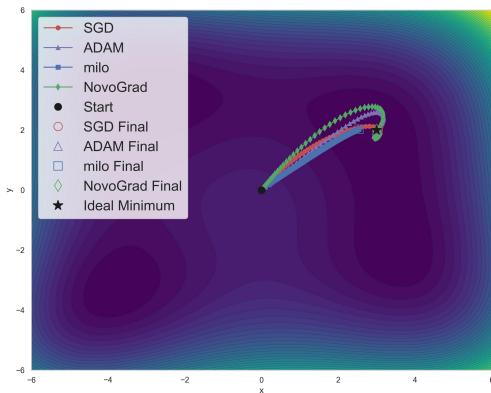
$$f(x_1, x_2) = (x_1^2 + x_2 - 11)^2 + (x_1 + x_2^2 - 7)^2. \quad (\text{D.1})$$

996 For the analytic gradient, we define:

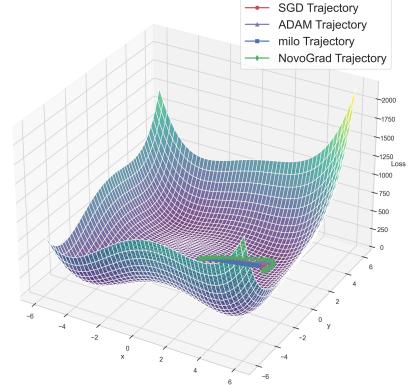
$$u = x_1^2 + x_2 - 11, \quad v = x_1 + x_2^2 - 7. \quad (\text{D.2})$$

997 Then, the partial derivatives are:

$$\frac{\partial f}{\partial x_1} = 4x_1 u + 2v, \quad \frac{\partial f}{\partial x_2} = 2u + 4x_2 v. \quad (\text{D.3})$$



(a) Himmelblau Loss 2D



(b) Himmelblau Loss 3D

Figure 6: Himmelblau Loss Function: 2D Contour vs. 3D Loss Map

998 Experimental Setup

- 999 • **Domain:** $[-6, 6] \times [-6, 6]$
- 1000 • **Initial Point:** $(0, 0)$
- 1001 • **Iterations:** 150

1002 Observations and Analysis

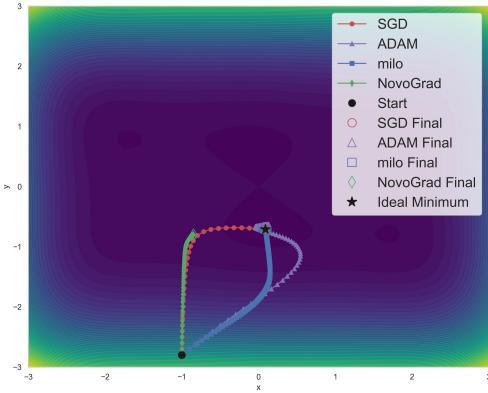
- 1003 • **SGD:** The SGD trajectory follows a generally direct path, quickly descending into a
 1004 low-loss valley, curving towards the end, and stabilizing near the minimum.
- 1005 • **Adam:** Adam has a less direct path, overshooting and eventually curving around to reach the objective.
- 1006 • **NovoGrad:** NovoGrad, despite taking normalized steps, overshoots with similar behavior to Adam,
 1007 eventually curving around to reach the objective.
- 1008 • **Milo:** Milo takes normalized steps, taking a more conservative path since the large gradients are scaled
 1009 down, following a uniform trajectory toward the basin.

- 1010 **Key Takeaway:** All four optimizers ultimately approach a global minimum, but SGD and Milo exhibit a more
 1011 stable and direct convergence compared to Adam and NovoGrad.

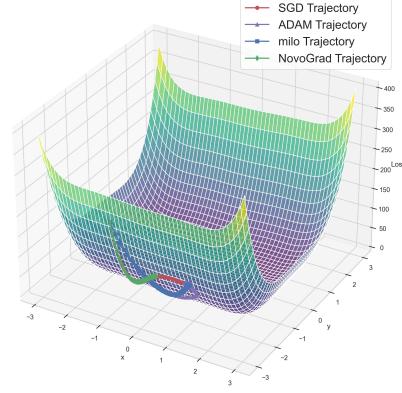
1012 **D.5.2 Camel6 Function**

1013 The Camel6 function is defined by:

$$f(x_1, x_2) = \left(4 - 2.1 x_1^2 + \frac{x_1^4}{3}\right) x_1^2 + x_1 x_2 + \left(-4 + 4 x_2^2\right) x_2^2. \quad (\text{D.4})$$



(a) Camel6 Loss 2D



(b) Camel6 Loss 3D

Figure 7: Camel6 Loss Function: 2D Contour vs. 3D Loss Map

1014 **Experimental Setup**

1015 • **Domain:** $[-3, 3] \times [-3, 3]$

1016 • **Initial Point:** $(-1, -2.8)$

1017 • **Iterations:** 150

1018 **Observations and Analysis**

1019 • **SGD:** The trajectory takes consistent steps, taking a largely indirect path to the local minima and
1020 eventually converging to the global.

1021 • **Adam:** Similar to SGD Adam takes an indirect path to the local minima, then shifts to a smaller step
1022 size with a more direct path to the eventual global minima.

1023 • **NovoGrad:** Novograd follows a similar path as SGD, however it comes to an early stop after it
1024 reaches the local minima, never reaching the global.

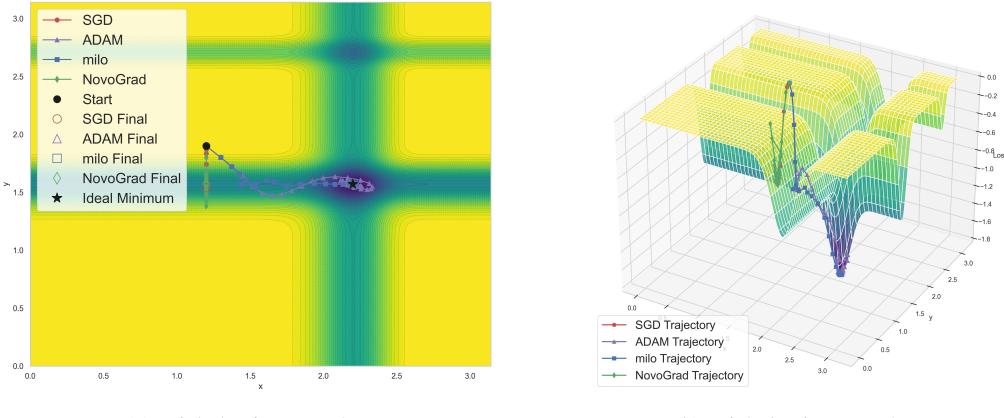
1025 • **Milo:** With normalized updates, Milo takes consistent steps, initially following the same path as Adam,
1026 but deviating towards the global minima for a more direct convergence than the other optimizers.

1027 **Key Takeaway:** In the Camel6 landscape, Milo's consistent step sizes help it maintain stable progress even
1028 when facing plateaus, with SGD and Adam taking more indirect paths. NovoGrad never reaches convergence in
1029 this problem, instead settling in the local minima it quickly reached.

1030 **D.5.3 Michalewicz Function**

1031 For a d -dimensional input x and a parameter m (typically $m = 10$), the Michalewicz function is:

$$f(x) = - \sum_{i=1}^d \sin(x_i) \left[\sin\left(\frac{i x_i^2}{\pi}\right) \right]^{2m}. \quad (\text{D.5})$$



(a) Michalewicz Loss 2D

(b) Michalewicz Loss 3D

Figure 8: Michalewicz Loss Function: 2D Contour vs. 3D Loss Map

1032 **Experimental Setup**

- 1033 • **Domain:** $[-10, 10] \times [-10, 10]$
- 1034 • **Initial Point:** $(2.4, -4.3)$
- 1035 • **Iterations:** 150

1036 **Observations and Analysis**

- 1037 • **SGD:** For the Michalewicz experiment, SGD quickly became stuck in the local minima within the basin, never reaching the global.
- 1038 • **Adam:** Adam's adaptive step size helps it overcome the local traps, reaching the global minima, however it takes a more indirect path and briefly overshoots the goal.
- 1039 • **NovoGrad:** NovoGrad, similar to SGD, also becomes quickly stuck in the local minima within the basin.
- 1040 • **Milo:** Milo successfully reaches the global minima, making consistent and direct progress towards the goal with no overshooting.

1041 **Key Takeaway:** De Jong's function, with its many local minima, highlights Milo's strength in sustaining movement when the gradient is very small. Adam adapts well but again overshoots the goal before circling back.
 1042 Additionally, it is shown that SGD and NovoGrad can become trapped in shallow basins.

1043 **D.5.4 De Jong Function 5**

1044 De Jong's fifth function is a multimodal test function defined for a 2-dimensional input vector $x = (x_1, x_2)$. It
 1045 is constructed using a fixed grid of constants and is known for its complex landscape with many local minima.
 1046 The function is given by:
 1047

$$f(x) = \left[0.002 + \sum_{i=1}^{25} \frac{1}{i + (x_1 - a_{1i})^6 + (x_2 - a_{2i})^6} \right]^{-1}, \quad (\text{D.6})$$

1048 where the constants a_{1i} and a_{2i} form a 5×5 grid of values from the set $\{-32, -16, 0, 16, 32\}$, specifically:

$$\begin{aligned} a_{1i} &= \text{tile}([-32, -16, 0, 16, 32], 5), \\ a_{2i} &= \text{repeat}([-32, -16, 0, 16, 32], 5). \end{aligned}$$

1049 This function is highly sensitive to small changes in x , making it useful for evaluating the performance of
 1050 optimization algorithms on rugged search landscapes.

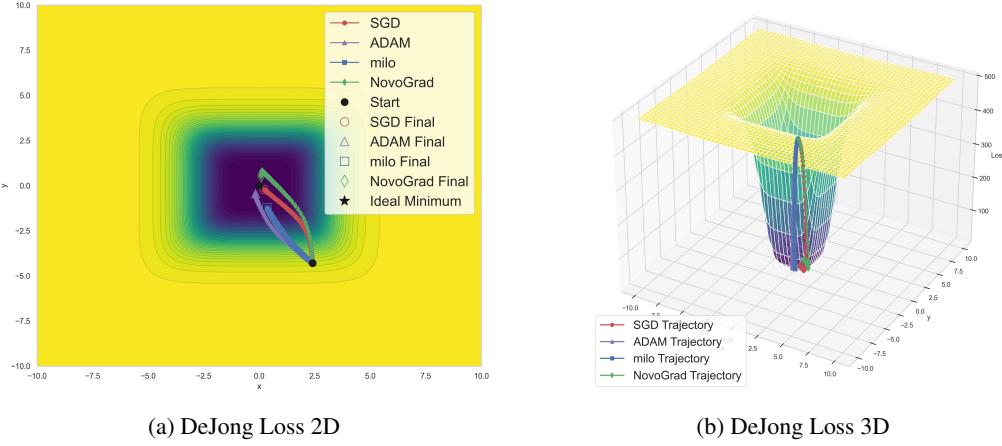


Figure 9: DeJong’s Fifth Function Loss Function: 2D Contour vs. 3D Loss Map

1055 Experimental Setup

1056 • **Domain:** $[0, \pi] \times [0, \pi]$

1057 • **Initial Point:** $(1.2, 1.9)$

1058 • **Iterations:** 250

1059 Observations and Analysis

1060 • **SGD:** Due to the steep descent structure of the function, SGD curves around the edge until it corrects
1061 towards the global minima, reaching it.

1062 • **Adam:** Adam’s adaptive updates help it follow a smoother descent, reaching the final global minima
1063 smoothly.

1064 • **NovoGrad:** NovoGrad follows a similar path around the edge as SGD, eventually overshooting
1065 slightly before reaching the global minima.

1066 • **Milo:** Milo’s normalized update rule maintains a consistent step size even when the gradients are
1067 small, allowing steady progress; however, it stops right as the bottom of the basin is reached, not
1068 continuing further as the other optimizers do.

1069 Key Takeaway:

1070 D.5.5 Experiment 1 Analysis

1071 The loss map experiments reveal that while SGD, Adam, NovoGrad, and Milo each have their strengths, the
1072 proposed Milo optimizer, with its normalized, group-based updates, can offer stable and consistent progress in
1073 challenging optimization landscapes. In particular, Milo shows promise in escaping flat regions and avoiding
1074 overshooting due to its consistent step size. However, as shown with DeJong, the consistent step size can
1075 potentially lead to it stopping earlier than desired.

Table 4: Key Experimental Settings, Experiments 2 - 4

Setting	Specification
Data splits	Validation: 15%, Test: 10%
Batch size	128
Epochs	10
Runs per optimizer	5
Hyperparameter tuning	5 trials per optimizer using log-uniform and uniform sampling over specified parameter grids
Base learning rates	LOGISTIC, MULTILAYER: 0.05; DEEPCNN, RESNET18: 0.005
Optimizers	Milo, Milo_LW, SGD, AdaGrad, AdamW, NovoGrad
Models & Datasets	LOGISTIC, MULTILAYER: MNIST; RESNET18: CIFAR-100
Transforms	MNIST: flatten to vector; CIFAR: tensor only (no augmentations)

1076 **D.6 Hyperparameter Tuning:**

Table 5: Hyperparameter search spaces for all optimizers

Optimizer	Search Space
SGD	lr: log-uniform [0.005, 0.1]; momentum: uniform [0.45, 0.99]; weight_decay: log-uniform [0.0005, 0.01]
ADAGRAD	lr: log-uniform [0.005, 0.1]; lr_decay: uniform [0.0, 0.1]; weight_decay: log-uniform [0.0005, 0.01]; eps: log-uniform [10^{-10} , 10^{-6}]
ADAMW	lr: log-uniform [0.005, 0.1]; weight_decay: log-uniform [0.0005, 0.01]; betas: {(0.9, 0.98), (0.95, 0.99)}; eps: log-uniform [10^{-9} , 10^{-6}]
NovoGrad	lr: log-uniform [0.005, 0.1]; betas: {(0.9, 0.98), (0.95, 0.99)}; weight_decay: log-uniform [0.0005, 0.01]
Milo	– (default settings)
Milo _{LW}	– (default settings)

Table 6: Default learning rates and optimizer parameter settings

Component	Default Setting
<i>Base learning rate</i>	
LOGISTIC	0.05
MULTILAYER	0.05
DEEPCNN	0.005
RESNET18	0.005
<i>Optimizer defaults</i>	
SGD	momentum=0.9; nesterov=True; weight_decay=0.0001
ADAGRAD	lr_decay=0; weight_decay=0.0; eps=1e-10
ADAMW	betas=(0.9, 0.999); eps=1e-8; weight_decay=0.01
Milo	layer_wise=False; scale_aware=True; scale_factor=0.2; momentum=0.9; use_cached_mapping=False; foreach=True
Milo _{LW}	layer_wise=True; scale_aware=True; scale_factor=0.2; momentum=0.9; use_cached_mapping=True; foreach=True
NovoGrad	betas=(0.9, 0.99); weight_decay=0.001; grad_averaging=True

1077 **D.7 Experiment 2: Logistic Regression**

1078 This experiment evaluates a convex multi-class logistic regression using the MNIST Dataset ([LeCun et al.,
 1079 1998]), just as was done in the study that proposed Adam. The goal is to compare optimizers with respect to
 1080 training loss, convergence speed, and test accuracy. The model uses L2 regularization with a single linear layer
 1081 mapping a 784-dimensional input to 10 classes. Visual results can be seen in Figure 3.

Table 7: Logistic Regression Model Architecture Layers

Layer	Description
Input Layer	Flattened input vector of dimension 784 (28×28)
Linear	Single fully connected layer mapping 784 to 10 (number of classes)

Table 8: Logistic Regression Training Results

Optimizer	Epoch Group	Loss	Accuracy (%)	F1 Score
Milo	1–4	0.2930	92.42	0.9233
	5–8	0.2732	93.07	0.9298
	9–10	0.2714	93.19	0.9310
Milo_{LW}	1–4	0.2594	92.87	0.9277
	5–8	0.2312	93.62	0.9354
	9–10	0.2258	93.75	0.9367
SGD	1–4	0.3076	91.47	0.9135
	5–8	0.2818	92.23	0.9212
	9–10	0.2773	92.34	0.9224
ADAGRAD	1–4	0.4866	88.03	0.8784
	5–8	0.4720	88.29	0.8810
	9–10	0.4687	88.33	0.8814
ADAMW	1–4	0.3103	91.48	0.9135
	5–8	0.2593	92.83	0.9273
	9–10	0.2490	93.18	0.9309
Novograd	1–4	0.3774	89.87	0.8973
	5–8	0.2928	91.83	0.9172
	9–10	0.2827	92.13	0.9202

Table 9: Average Runtime and Memory Increase for Logistic Regression Experiment (5 runs)

Optimizer	Avg Runtime (s)	Avg Memory (MB)
MILO	32.51	63.5
MILO_{LW}	32.66	2.0
SGD	29.71	0.0
ADAGRAD	29.76	0.0
ADAMW	30.03	16.4
NOVOGRAD	32.16	4.1

1082 **D.7.1 Statistical Analysis**

Table 10: Pairwise Significance Tests for Validation Loss

Optimizer A	Optimizer B	Mean A	Mean B	Better	p-value	Sig.
MILO	MILO _{LW}	0.350833	0.304008	MILO _{LW}	8.13499e-08	***
MILO	SGD	0.350833	0.306714	SGD	3.62048e-06	***
MILO	ADAMW	0.350833	0.287661	ADAMW	9.31980e-07	***
MILO	ADAGRAD	0.350833	0.481010	MILO	9.90140e-11	***
MILO	NOVOGRAD	0.350833	0.309194	NOVOGRAD	4.23780e-06	***
MILO _{LW}	SGD	0.304008	0.306714	MILO _{LW}	1.32372e-01	
MILO _{LW}	ADAMW	0.304008	0.287661	ADAMW	9.06983e-05	***
MILO _{LW}	ADAGRAD	0.304008	0.481010	MILO _{LW}	2.49227e-13	***
MILO _{LW}	NOVOGRAD	0.304008	0.309194	MILO _{LW}	1.67816e-02	*
SGD	ADAMW	0.306714	0.287661	ADAMW	1.66218e-08	***
SGD	ADAGRAD	0.306714	0.481010	SGD	2.60849e-11	***
SGD	NOVOGRAD	0.306714	0.309194	SGD	2.46024e-02	*
ADAMW	ADAGRAD	0.287661	0.481010	ADAMW	4.96532e-11	***
ADAMW	NOVOGRAD	0.287661	0.309194	ADAMW	9.13035e-09	***
ADAGRAD	NOVOGRAD	0.481010	0.309194	NOVOGRAD	1.79728e-11	***

Table 11: Pairwise Significance Tests for Validation F1-Score

Optimizer A	Optimizer B	Mean A	Mean B	Better	p-value	Sig.
MILO	MILO _{LW}	0.912861	0.916955	MILO _{LW}	1.75279e-03	**
MILO	SGD	0.912861	0.914048	SGD	2.41945e-01	
MILO	ADAGRAD	0.912861	0.876378	MILO	5.42161e-10	***
MILO	ADAMW	0.912861	0.919612	ADAMW	5.82021e-05	***
MILO	NOVOGRAD	0.912861	0.913075	NOVOGRAD	7.63541e-01	
MILO _{LW}	SGD	0.916955	0.914048	MILO _{LW}	1.56224e-02	*
MILO _{LW}	ADAGRAD	0.916955	0.876378	MILO _{LW}	2.24596e-10	***
MILO _{LW}	ADAMW	0.916955	0.919612	ADAMW	9.94675e-03	**
MILO _{LW}	NOVOGRAD	0.916955	0.913075	MILO _{LW}	1.93199e-03	**
SGD	ADAGRAD	0.914048	0.876378	SGD	3.79011e-10	***
SGD	ADAMW	0.914048	0.919612	ADAMW	4.33298e-04	***
SGD	NOVOGRAD	0.914048	0.913075	SGD	2.51483e-01	
ADAGRAD	ADAMW	0.876378	0.919612	ADAMW	8.47333e-09	***
ADAGRAD	NOVOGRAD	0.876378	0.913075	NOVOGRAD	9.06177e-08	***
ADAMW	NOVOGRAD	0.919612	0.913075	ADAMW	3.71358e-05	***

1083 **D.8 Experiment 3: Multilayer Neural Networks**

1084 The objective of this experiment is to test the optimizer on a non-convex problem by using a multilayer perception.
 1085 For this study, our model choice was driven by the previous Adam publication, using two fully connected layers
 1086 with 1000 ReLU units each, with dropout and weight decay on a minibatch of 128, on the MNIST image dataset
 1087 ([LeCun et al., 1998]).

Table 12: Multilayer NN Model Architecture Layers

Layer	Description
Linear	Fully connected layer: input dimension 784 to hidden dimension 128
ReLU	Activation function
Dropout	Dropout layer with probability 0.2
Linear	Fully connected layer: 128 to 10 (number of classes)

Table 13: Multilayer Network Training Results

Optimizer	Epoch Group	Loss	Accuracy (%)	F1 Score
Milo	1–4	0.0777	98.24	0.9822
	5–8	0.0161	99.66	0.9966
	9–10	0.0050	99.90	0.9990
Milo_{LW}	1–4	0.1034	96.95	0.9692
	5–8	0.0372	98.88	0.9887
	9–10	0.0215	99.36	0.9936
SGD	1–4	0.1695	95.32	0.9528
	5–8	0.0921	97.63	0.9762
	9–10	0.0756	98.16	0.9815
ADAGRAD	1–4	0.2204	93.98	0.9391
	5–8	0.2013	94.53	0.9447
	9–10	0.1968	94.66	0.9461
ADAMW	1–4	0.0791	97.58	0.9756
	5–8	0.0259	99.14	0.9913
	9–10	0.0144	99.51	0.9951
Novograd	1–4	0.1655	95.24	0.9519
	5–8	0.0637	98.18	0.9817
	9–10	0.0461	98.72	0.9871

Table 14: Average Runtime and Memory Increase for Multilayer Experiment (5 runs)

Optimizer	Avg Duration (s)	Avg Memory (MB)
MILO	34.37	18.4
MILO_{LW}	36.05	0.0
SGD	30.28	0.0
ADAGRAD	30.74	0.0
ADAMW	30.87	0.0
NOVOGRAD	34.24	0.0

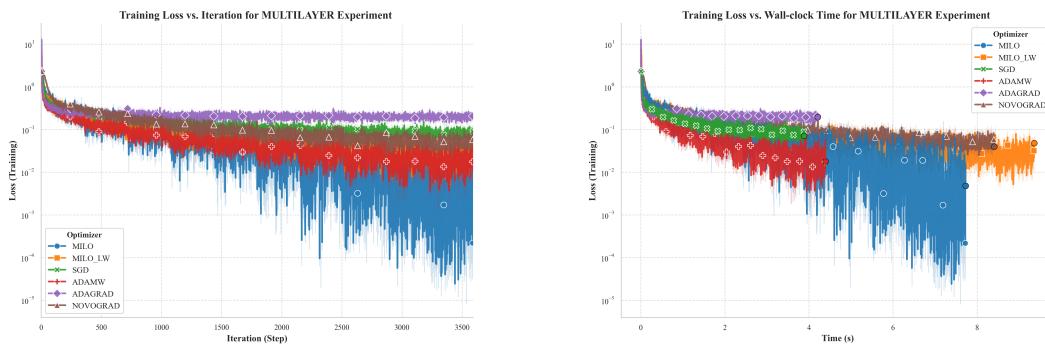


Figure 10: Comparison of MLP training cost by iterations and wall time.

1088 **D.8.1 Statistical Analysis**

Table 15: Pairwise Significance Tests for Validation Loss (Final Epoch)

Optimizer A	Optimizer B	Mean A	Mean B	Better	p-value	Sig.
MILO	MILO _{LW}	0.225943	0.150747	MILO _{LW}	1.51986e-04	***
MILO	SGD	0.225943	0.104657	SGD	1.10944e-04	***
MILO	ADAMW	0.225943	0.121495	ADAMW	2.48203e-05	***
MILO	ADAGRAD	0.225943	0.215592	ADAGRAD	3.16927e-01	
MILO	NOVOGRAD	0.225943	0.0872035	NOVOGRAD	2.96707e-05	***
MILO _{LW}	SGD	0.150747	0.104657	SGD	4.70080e-04	***
MILO _{LW}	ADAMW	0.150747	0.121495	ADAMW	1.64005e-03	**
MILO _{LW}	ADAGRAD	0.150747	0.215592	MILO _{LW}	1.31936e-05	***
MILO _{LW}	NOVOGRAD	0.150747	0.0872035	NOVOGRAD	2.06945e-05	***
SGD	ADAMW	0.104657	0.121495	SGD	1.71513e-02	*
SGD	ADAGRAD	0.104657	0.215592	SGD	2.16145e-05	***
SGD	NOVOGRAD	0.104657	0.0872035	NOVOGRAD	9.02238e-04	***
ADAMW	ADAGRAD	0.121495	0.215592	ADAMW	7.75432e-07	***
ADAMW	NOVOGRAD	0.121495	0.0872035	NOVOGRAD	4.28650e-04	***
ADAGRAD	NOVOGRAD	0.215592	0.0872035	NOVOGRAD	1.16408e-06	***

Table 16: Pairwise Significance Tests for Validation F1-Score (Final Epoch)

Optimizer A	Optimizer B	Mean A	Mean B	Better	p-value	Sig.
MILO	MILO _{LW}	0.973116	0.968024	MILO	9.10593e-04	***
MILO	SGD	0.973116	0.971065	MILO	5.32615e-02	
MILO	ADAMW	0.973116	0.972800	MILO	8.14991e-01	
MILO	ADAGRAD	0.973116	0.941102	MILO	6.55427e-06	***
MILO	NOVOGRAD	0.973116	0.973147	NOVOGRAD	9.80252e-01	
MILO _{LW}	SGD	0.968024	0.971065	SGD	1.56224e-02	**
MILO _{LW}	ADAMW	0.968024	0.972800	ADAMW	6.94322e-03	**
MILO _{LW}	ADAGRAD	0.968024	0.941102	MILO _{LW}	3.03333e-05	***
MILO _{LW}	NOVOGRAD	0.968024	0.973147	NOVOGRAD	2.48226e-03	**
SGD	ADAMW	0.971065	0.972800	ADAMW	1.83324e-01	
SGD	ADAGRAD	0.971065	0.941102	SGD	3.85788e-05	***
SGD	NOVOGRAD	0.971065	0.973147	NOVOGRAD	8.87216e-02	
ADAMW	ADAGRAD	0.972800	0.941102	ADAMW	1.98429e-06	***
ADAMW	NOVOGRAD	0.972800	0.973147	NOVOGRAD	8.11459e-01	
ADAGRAD	NOVOGRAD	0.941102	0.973147	NOVOGRAD	3.13313e-06	***

1089 **D.9 Experiment 4: Convolutional Neural Networks (CNNs)**

1090 This experiment evaluates a deep Residual Network, specifically ResNet-18 [He et al., 2016], on the CIFAR-100
 1091 dataset [Krizhevsky and Hinton, 2009] to test optimizer performance in a finer-grained classification setting.
 1092 Compared to CIFAR-10, CIFAR-100 has 100 target classes, which stresses both the representation capacity of
 1093 the network and the optimizer’s ability to generalize. We use the standard ResNet-18 architecture, modifying
 1094 only the final linear layer to output 100 logits. Input images are augmented with 4-pixel zero-padding, then
 1095 per-channel normalized to zero mean and unit variance. All training runs use a batch size of 128.

Table 17: Deep CNN Model Architecture Layers

Layer	Description
Conv2d	Input: 3 channels, Output: 32 channels, Kernel size: 3, Padding: 1
ReLU	Activation function
Conv2d	Input: 32 channels, Output: 32 channels, Kernel size: 3, Padding: 1
ReLU	Activation function
MaxPool2d	Pooling layer with kernel size 2
Conv2d	Input: 32 channels, Output: 64 channels, Kernel size: 3, Padding: 1
ReLU	Activation function
Conv2d	Input: 64 channels, Output: 64 channels, Kernel size: 3, Padding: 1
ReLU	Activation function
MaxPool2d	Pooling layer with kernel size 2
Flatten	Converts feature maps into a 1D vector
Linear	Fully connected layer: 4096 (i.e. $64 \times 8 \times 8$) to 128
ReLU	Activation function
Linear	Fully connected layer: 128 to 10 (number of classes)

Table 18: ResNet-18 Training Results

Optimizer	Epoch Group	Loss	Accuracy (%)	F1 Score
Milo	1–4	2.6884	33.03	0.3130
	5–8	1.4598	61.59	0.6113
	9–10	0.7754	80.63	0.8064
Milo_{LW}	1–4	3.0691	23.40	0.2103
	5–8	2.0448	44.88	0.4389
	9–10	1.5241	57.17	0.5651
SGD	1–4	2.6914	34.53	0.3335
	5–8	1.2595	63.90	0.6395
	9–10	0.6575	81.21	0.8144
ADAGRAD	1–4	2.7093	32.94	0.3112
	5–8	1.7878	55.53	0.5470
	9–10	1.4423	65.83	0.6533
ADAMW	1–4	3.2084	22.69	0.2104
	5–8	1.1366	68.65	0.6847
	9–10	0.2509	92.53	0.9258
NOVOGRAD	1–4	3.1196	24.28	0.2205
	5–8	1.1359	66.98	0.6673
	9–10	0.3902	87.82	0.8780

Table 19: Average Runtime and Memory Increase for ResNet-18 Experiment (5 runs)

Optimizer	Avg Duration (s)	Avg Memory (MB)
MILO	148.53	59.4
MILO_{LW}	182.82	2.0
SGD	130.56	14.3
ADAGRAD	130.69	0.0
ADAMW	131.10	14.3
NOVOGRAD	172.25	4.1

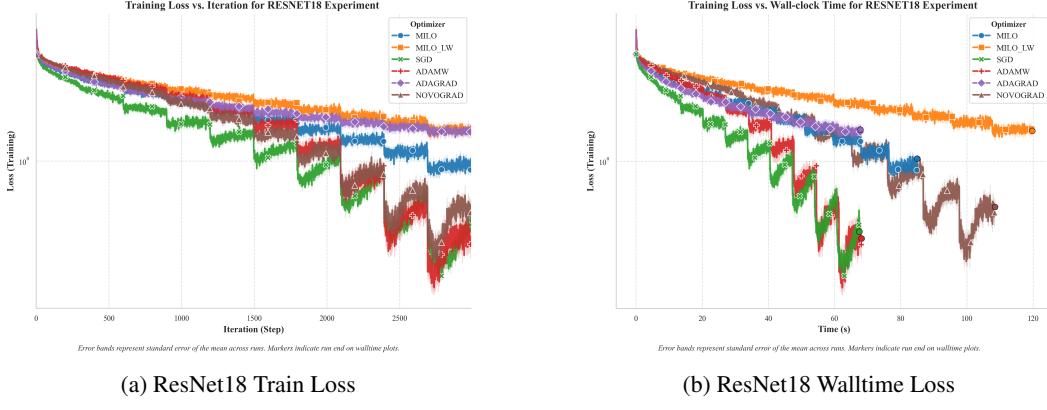


Figure 11: Comparison of ResNet18 training cost by iterations and wall time.

1096 D.9.1 Statistical Analysis

Table 20: Pairwise Significance Tests for Validation Loss

Opt. A	Opt. B	Mean A	Mean B	Better	p-value	Sig.
MILO	MILO _{LW}	2.09638	2.11451	MILO	0.428833	
MILO	SGD	2.09638	2.16893	MILO	0.044138	*
MILO	ADAGRAD	2.09638	2.48136	MILO	4.41550e-07	***
MILO	ADAMW	2.09638	3.38785	MILO	1.17780e-06	***
MILO	NOVOGRAD	2.09638	2.29087	MILO	0.000535	***
MILO _{LW}	SGD	2.11451	2.16893	MILO _{LW}	0.081725	
MILO _{LW}	ADAGRAD	2.11451	2.48136	MILO _{LW}	1.24390e-06	***
MILO _{LW}	ADAMW	2.11451	3.38785	MILO _{LW}	6.11940e-06	***
MILO _{LW}	NOVOGRAD	2.11451	2.29087	MILO _{LW}	0.001571	**
SGD	ADAGRAD	2.16893	2.48136	SGD	9.18010e-06	***
SGD	ADAMW	2.16893	3.38785	SGD	6.27030e-07	***
SGD	NOVOGRAD	2.16893	2.29087	SGD	0.009368	**
ADAGRAD	ADAMW	2.48136	3.38785	ADAGRAD	8.92350e-06	***
ADAGRAD	NOVOGRAD	2.48136	2.29087	NOVOGRAD	0.000613	***
ADAMW	NOVOGRAD	3.38785	2.29087	NOVOGRAD	6.09120e-07	***

Table 21: Pairwise Significance Tests for Validation F1-Score

Opt. A	Opt. B	Mean A	Mean B	Better	p-value	Sig.
MILO	MILO _{LW}	0.481343	0.441000	MILO	0.000321	***
MILO	SGD	0.481343	0.454009	MILO	0.003989	**
MILO	ADAGRAD	0.481343	0.365256	MILO	1.58180e-07	***
MILO	ADAMW	0.481343	0.411299	MILO	0.000352	***
MILO	NOVOGRAD	0.481343	0.492262	NOVOGRAD	0.113014	
MILO _{LW}	SGD	0.441000	0.454009	SGD	0.050683	
MILO _{LW}	ADAGRAD	0.441000	0.365256	MILO _{LW}	1.69300e-06	***
MILO _{LW}	ADAMW	0.441000	0.411299	MILO _{LW}	0.024678	*
MILO _{LW}	NOVOGRAD	0.441000	0.492262	NOVOGRAD	3.15430e-06	***
SGD	ADAGRAD	0.454009	0.365256	SGD	6.67760e-07	***
SGD	ADAMW	0.454009	0.411299	SGD	0.005071	**
SGD	NOVOGRAD	0.454009	0.492262	NOVOGRAD	0.000213	***
ADAGRAD	ADAMW	0.365256	0.411299	ADAMW	0.003599	**
ADAGRAD	NOVOGRAD	0.365256	0.492262	NOVOGRAD	4.01170e-08	***
ADAMW	NOVOGRAD	0.411299	0.492262	NOVOGRAD	0.000330	***

1097 **D.10 Experiment 5: Reinforcement Learning**

1098 Reinforcement Learning (RL) has become increasingly vital across various domains in recent years. RL enables
 1099 agents to learn optimal behavior across a variety of environments through trial and error. This results in agents
 1100 being more adaptable than traditional methods in complex environments. It plays a role in autonomous systems,
 1101 sequential decision-making without supervision, and LLMS through Reinforcement Learning from Human
 1102 Feedback (RLHF). Due to these growing uses, this study explores the performance of our optimizer in a three-
 1103 dimensional space with stable consistent points, and randomly generated points. The purpose of this experiment
 1104 is to study the optimizer's influence on reward-based actions and its adaptability. This experiment involved
 1105 simulating a boxed environment, providing the agent with the observation of its current location and goal, and
 1106 simulating 1000 episodes with 500 steps each. Each agent was trained 5 individual times, given 5 runs per
 1107 training regimen, resulting in the final results being the averaged prediction of each agent across 25 runs, with
 1108 the success rate being the rate the agent correctly reached the goal. The following results were observed.

Table 22: RL Experiment Training Parameters

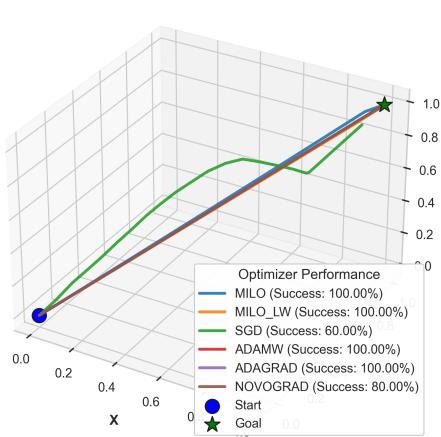
Parameter	Value
Learning Rate	0.0005
Gamma (Discount Factor)	0.99
Episodes	1000
Log Interval	250 episodes
Initial Exploration	0.2
Final Exploration	0.01
Exploration Decay	0.998
Gradient Clip	0.5
Distance Threshold	0.05
Random Goal	True
Max Steps per Episode	500
Start Position	[0.0, 0.0, 0.0]
Goal Position	[1, 1, 1]

Table 23: Policy Network Architecture

Component	Description
Input Embedding	A linear layer mapping the input (observation) to 256 dimensions, followed by LayerNorm and ReLU activation.
Residual Blocks	4 residual blocks; each block consists of: <ul style="list-style-type: none"> • Linear ($256 \rightarrow 128$) • LayerNorm + ReLU • Dropout (0.1) • Linear ($128 \rightarrow 256$) • LayerNorm • Skip connection
Output Layers	A sequential bottleneck: Linear ($256 \rightarrow 128$) + LayerNorm + ReLU, Linear ($128 \rightarrow 64$) + LayerNorm + ReLU, and Linear ($64 \rightarrow$ output dimension), with final outputs scaled via \tanh (multiplied by 0.1).

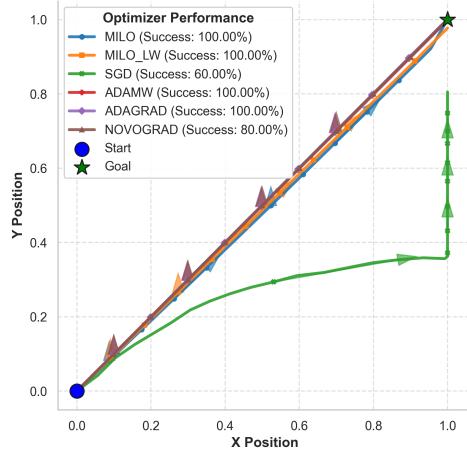
1109 **D.10.1 Set Result:**

Average Agent Trajectories by Optimizer



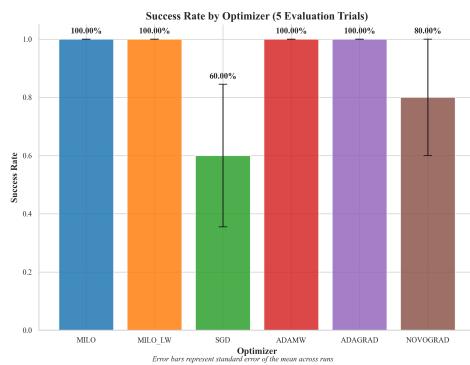
(a) 3D Agent Trajectories

Top-Down View of Average Agent Trajectories



(b) 2D Agent Trajectories

Figure 12: 3D and 2D visuals of set RL agent trajectories.



(a) RL Set Success Rates



(b) Reward Comparison

Figure 13: Comparison of set success rates by optimizer and average reward by episode.

1110 As shown in Figure 12 and Figure 13, each of the optimizers has a high success rate of reaching the set goal
 1111 with no random initialization or generation of the goal. The optimizers with the highest performance proved
 1112 to be AdamW, AdaGrad, Milo, and Milo_{LW}, each achieving a 100% accuracy. Overall, the set experiment
 1113 demonstrates that each of the models can effectively learn the patterns generated in this 3D space.

1114 **D.10.2 Randomized Result:**

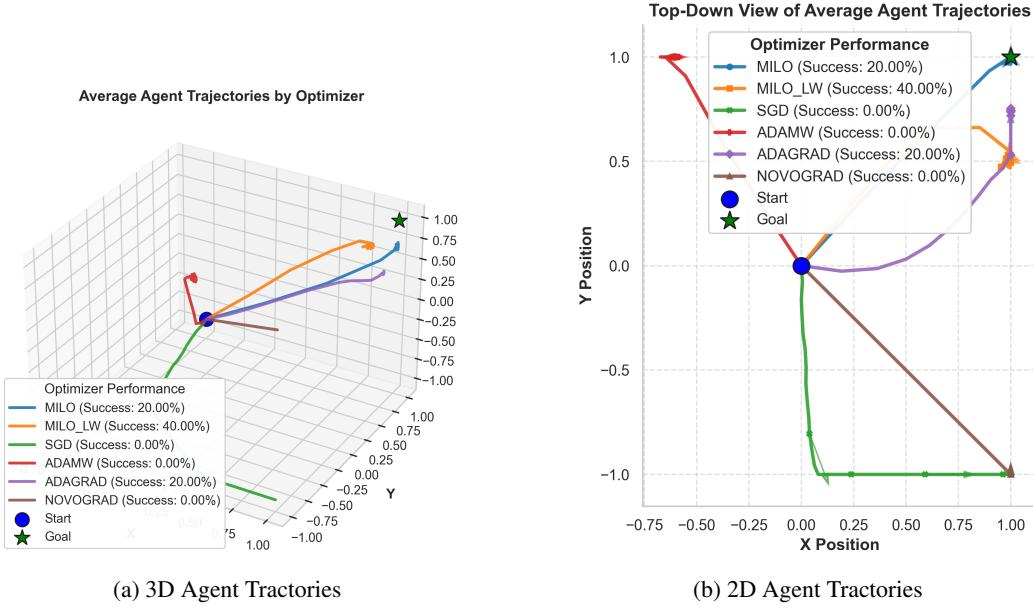


Figure 14: 3D and 2D visuals of set RL agent trajectories.

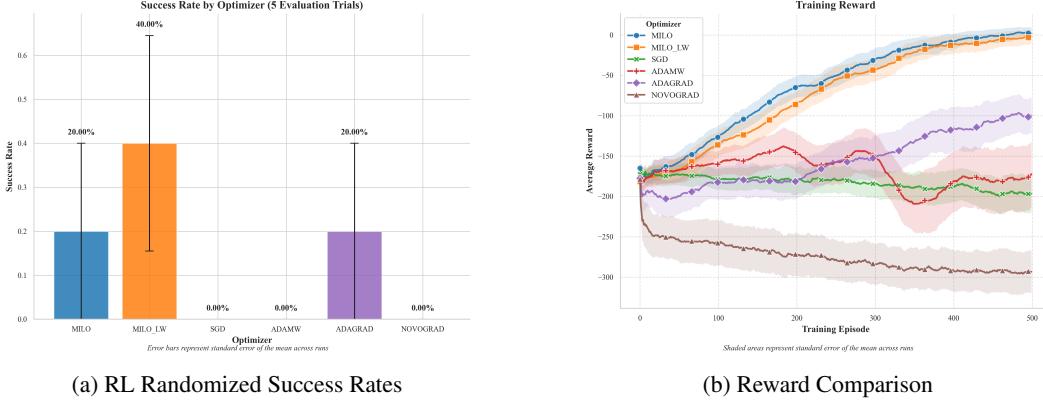


Figure 15: Comparison of randomized success rates by optimizer and average reward by episode.

1115 The results of the randomly generated start and goal differ drastically from the set, with each optimizer performing
 1116 worse than before. However, each optimizer retains the relative performance shown previously regarding their
 1117 performance against one another. As shown in Figure 14a, Figure 14b, and Figure 15a, Milo_{LW} has the highest
 1118 accuracy, with a final success rate of 40%, being a large improvement over the success rates of the other
 1119 optimizers. This performance is also shown in Figure 15b, with Milo_{LW} demonstrating a higher reward than the
 1120 other optimizers.

1121 **D.11 Experiment 6: Large Language Model**

1122 LLMs have experienced significant growth in popularity in recent years, accompanied by numerous innovations
 1123 in model architecture and optimization. This study explores the performance of our optimizer in a pretraining
 1124 setting, using a small-scale Transformer model trained from scratch on a processed subset of The Pile dataset
 1125 [Gao et al., 2020]. The model weights are randomly initialized at the start of training, and no pre-trained
 1126 parameters are used at any stage. This setup allows us to evaluate optimizer behavior in a controlled language
 1127 modeling task, where both model capacity and training steps are deliberately constrained to enable rapid
 1128 comparison across runs.

1129 To support this investigation, we adapt the open-source training framework of FareedKhan-dev (2025), which
 1130 implements a simplified version of the Transformer architecture introduced by Vaswani et al. (2017). This
 1131 architecture includes core components such as multi-head self-attention, layer normalization, and position-wise
 1132 feed-forward networks, implemented natively in PyTorch. We extend the original implementation to support
 1133 multiple optimizers and perform controlled training and evaluation across several runs. Visual results can be
 1134 seen in Figure 4.

Table 24: Training Parameters

Parameter	Value	Description
VOCAB_SIZE	50304	Number of unique tokens in the vocabulary.
CONTEXT_LENGTH	64	Maximum sequence length for the model.
N_EMBED	64	Dimension of the embedding space.
N_HEAD	4	Number of attention heads in each transformer block.
N_BLOCKS	1	Number of transformer blocks in the model.
T_BATCH_SIZE	32	Number of samples per training batch.
T_CONTEXT_LENGTH	16	Context length for training batches.
T_TRAIN_STEPS	3000	Total number of training steps.
T_EVAL_STEPS	500	Frequency (in steps) for evaluation.
T_EVAL_ITERS	125	Number of iterations for evaluation.
T_LR_DECAY_STEP	1500	Step at which the learning rate decays.
T_LR	0.05	Initial learning rate for training.
T_LR_DECAYED	0.005	Learning rate after decay.
DEVICE	cuda	Device used for training.

Table 25: Model Architecture Layers

Layer	Description
Token Embedding	Converts input token indices into embeddings of dimension N_EMBED .
Positional Embedding	Adds positional information to the token embeddings for a maximum sequence length of CONTEXT_LENGTH.
Transformer Block	Comprises a multi-head self-attention mechanism (with N_HEAD heads) and a feed-forward network, incorporating residual connections and layer normalization.
Final Linear Projection	Projects the transformer output to logits over the vocabulary (transforming from N_EMBED to VOCAB_SIZE).

Table 26: Transformer Experiment Configuration

Component	Details
Model	Transformer
Vocabulary Size	50304
Embedding Dimension	64
Attention Heads	4
Transformer Blocks	1
Training Batch Size	32
Training Context Length	16
Total Training Steps	3000
Evaluation Frequency	Every 500 steps (125 iterations)
Learning Rate	Initial: 0.05, decays to 0.005 at step 1500

1135 The results in this experiment, as shown in Figure 4a and Figure 4b, demonstrate that Milo and Milo_{LW}
 1136 demonstrate comparable performance to some of the industry standard used for training LLMs such as Adam
 1137 and AdamW, with Milo_{LW} and Milo achieving the lowest recorded loss and fastest reduction of loss as well.
 1138 Additionally, since this was a more complex task, the wall time of each optimizer was logged as well, with Milo
 1139 being comparable to Adam and AdamW, whereas Milo_{LW} and NovoGrad took a longer training time.

1140 **D.12 Experiment 7: Ablation Study**

1141 To assess the impact of Milo’s design choices, we conduct a comprehensive ablation study across two model
 1142 architectures (MLP, ResNet-18) and two datasets (MNIST, CIFAR-10). Our goals are to isolate the effects of
 1143 key hyperparameters—group size, learning rate, clipping norm, and scale-aware blending—and evaluate their
 1144 impact on training stability, convergence speed, and sensitivity to hyperparameter selection.

1145 **Experimental Setup.** We evaluate two optimizer variants: standard Milo and its layer-wise variant Milo_{LW}.
 1146 For each model–dataset pair, we establish a baseline configuration and systematically vary one hyperparameter
 1147 at a time. The key parameters we will discuss include:

- 1148 • **Learning Rate (lr):** {0.0001, 0.001, 0.005, 0.01}
- 1149 • **Momentum (momentum):** {0.0, 0.5, 0.9}
- 1150 • **Weight Decay (weight_decay):** {0, 1e-4, 1e-3}
- 1151 • **Epsilon (eps):** {1e-8, 1e-5, 1e-2}
- 1152 • **Scale-Aware Blending (scale_aware, scale_factor):** {True, False}, {0.1, 0.2, 0.5} (Milo_{LW} only)

1153 We also explore the other parameters used by the model in our code, but only discuss these here as they had the
 1154 most notable results.

1155 We measure:

- 1156 • *Convergence Speed:* epochs to reach fixed loss threshold.
- 1157 • *Stability:* variance in final loss and accuracy across seeds.
- 1158 • *Hyperparameter Sensitivity:* performance variation relative to baseline.

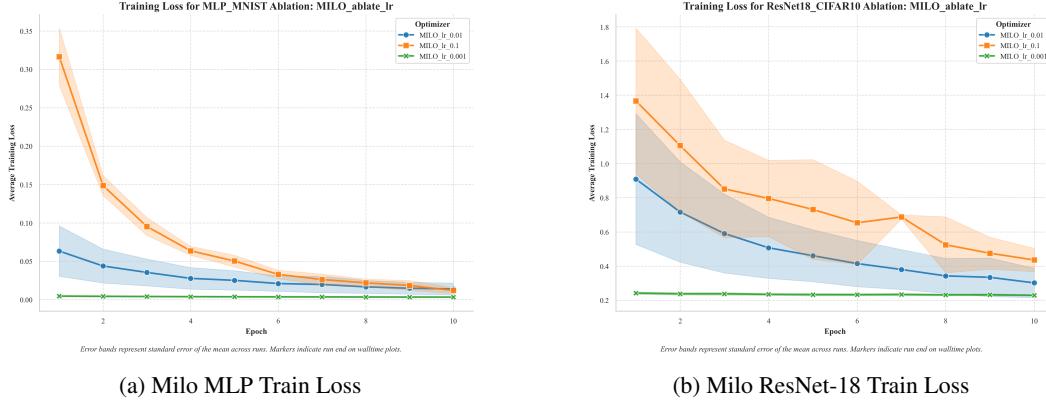


Figure 16: Comparison of Milo LR Ablation in MLP and ResNet-18.

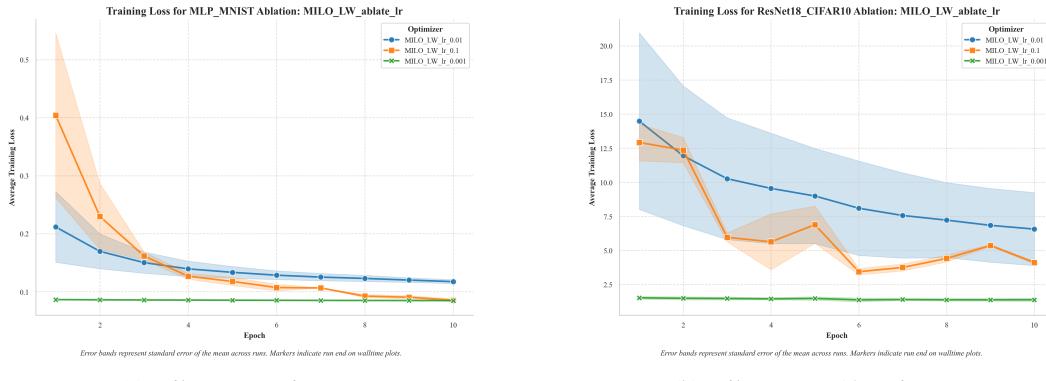


Figure 17: Comparison of Milo_{LW} LR Ablation in MLP and ResNet-18.

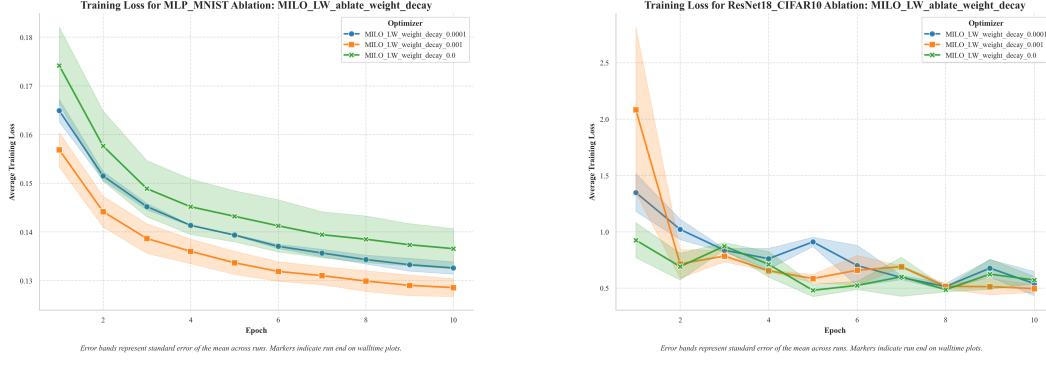


Figure 18: Comparison of Milo_LW Weight Decay Ablation in MLP and ResNet-18.

1159 **Key Findings.** We found that **learning rate** remains the most critical factor: optimal values vary by task and
 1160 model, with Milo and Milo_{LW} exhibiting similar sensitivity profiles. **Momentum** and **scale-aware blending**
 1161 have a moderate impact, influencing the convergence rate and final plateauing of the model. However, there is not
 1162 as large of a deviation in accuracy as changes in the learning rate. Lastly, **weight decay**, **epsilon**, and a majority
 1163 of the remaining hyperparameters were found to have minimal affect on the training. These results suggest that
 1164 Milo's group normalization and adaptive components are robust across a wide range of hyperparameters, with
 1165 the LR being the primary parameter needing to be tuned. This reduces the tuning burden relative to conventional
 1166 optimizers, as discussed earlier in this paper.