

# 385 Appendices

## 386 Contents

387	<b>1 Introduction</b>	1
388	<b>2 Related Work</b>	2
389	<b>3 Algorithm</b>	3
390	<b>4 Convergence Guarantees for Milo</b>	6
391	<b>5 Experiments</b>	7
392	<b>6 Limitations and Future Work</b>	9
393	<b>7 Conclusion</b>	9
394	<b>Appendices</b>	13
395	<b>A Nomenclature</b>	14
396	<b>B Convergence Proof for Convex Settings</b>	15
397	<b>C Convergence Proof for Non-Convex Settings</b>	20
398	<b>D Experimentation</b>	23

Table 2: Nomenclature

Symbol	Description
<b>Model Parameters and Tensors</b>	
$p_i$	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.
<b>Gradients and Grouping</b>	
$g_i$	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 $G_j$ .
$\mu_j$	Empirical mean of the components in group $G_j$ .
$\sigma_j^2$	Empirical variance of the components in group $G_j$ .
$\epsilon$	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 $i \in G_j$ .
$\rho_j$	Variance contraction factor for group $G_j$ , defined as $\rho_j = \frac{\sigma_j^2}{(\sigma_j + \epsilon)^2}$ .
$\rho$	Global constant such that $\rho_j \leq \rho < 1$ for all groups.
<b>Layer-Wise Grouping (if applicable)</b>	
$\mathcal{G}_l$	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 $l$ .
<b>Optimization Variables and Algorithm Parameters</b>	
$w_t$	Current model parameters (iterate) at iteration $t$ .
$f$ or $F$	Objective function (convex: $f$ ; non-convex: $F$ ).
$\eta$	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.
<b>Theoretical Analysis Parameters</b>	
$\delta_t$	Alignment error vector defined by $\tilde{g}_t = \nabla f(w_t) + \delta_t$ .
$\varepsilon$	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 $\varepsilon = \frac{\epsilon}{\sigma_{\min}} C(\tau)$ .
$\alpha$	Constant in $(0, 1)$ controlling the ratio between $\epsilon$ and $\sigma_j$ to ensure proper scaling.
$C(\tau)$	Constant that depends on the sub-Gaussian parameter $\tau$ of the gradient components.
$\Theta(\cdot)$	Tight asymptotic bound.
$O(\cdot)$	Big-O notation.
<b>Regret Analysis</b>	
$R(T)$	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).

400 **B Convergence Proof for Convex Settings**

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

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

404 Under standard  $L$ -smoothness, bounded-gradient, and variance-contraction/moment conditions, we obtain the  
405 following.

406 *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

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

407 showing improved dependence on the variance.

408

409 **B.2 Definitions and Assumptions**

410 **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  $\lambda \in [0, 1]$ ,

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

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

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

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

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

415 *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  
416  $L_{\max} = \max\{L_1, \dots, L_n\}$ .*

417 [See, e.g., [\[Nesterov 1983\]](#) for the smoothness condition.]

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

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

419 *with  $G > 0$ .*

420 [This is standard in stochastic optimization; see, e.g., [\[Bottou et al. 2018\]](#).]

421 **Assumption B.3** (Variance Contraction via Normalization). *Let  $g_t = \nabla f(w_t)$  be the gradient at iteration  $t$ .  
422 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  
423 mean and variance as*

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

424 *The normalized gradient is then defined by*

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

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

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

427 *Define the variance contraction factor as*

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

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

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

429 [See, e.g., [\[Wu and He 2018\]](#) for the effects of group normalization on variance contraction.]

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

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

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

434 *concentrate around the true moments.*

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

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

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

437 [Moment conditions as in [\[Vershynin 2018\]](#).]

438 **B.3 Auxiliary Lemmas**

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

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

441 A similar bound holds for the empirical variance.

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

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

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

445 where  $z_i$  are independent sub-Gaussian random variables with zero mean, unit variance, and sub-Gaussian

446 parameter  $\tau$ . Define the normalized elements by

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

447 Then,

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

448 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

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

449 Now, writing the decomposition

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

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

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

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

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

453 **Lemma B.3** (Robustness for Small  $\sigma_j$  with Explicit  $\epsilon$ ). *Assume that for each group  $G_j$  either:*

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

466 where  $\delta_t$  represents the alignment error (as introduced in Lemma B.2).

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

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

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

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

469 Here:

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

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

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

477 Summing over all groups, we obtain

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

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

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

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

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

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

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

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

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

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

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

485 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

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

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

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

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

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

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

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

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

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

490 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

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

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

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

492 thus demonstrating the advantage of variance reduction.

## 493 B.5 Discussion

494 The final convergence bound depends explicitly on the alignment error  $\varepsilon$ , the variance contraction factor  $\rho$ , and  
495 the learning rate  $\eta$ :

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

496 Lemmas B.1 and B.2 ensure that with sufficiently large group sizes and proper moment conditions (Assumption B.4), the empirical estimates concentrate and the normalized gradient remains well aligned with the true  
497 gradient. By choosing  $\epsilon$  relative to the minimum variance (as in Lemma B.3), even cases with very small  $\sigma_j$  are  
498 robustly handled. This completes the proof of the  $O(1/\sqrt{T})$  convergence rate and confirms that Milo yields  
499 improved regret bounds compared to standard SGD.  
500

## 501 B.6 Conclusion

502 We have presented a detailed, self-contained convergence proof for Milo under convex settings. By introducing  
503 the necessary definitions, assumptions, and auxiliary lemmas inline, and then integrating them into the conver-  
504 gence analysis, we demonstrate that under appropriate conditions the alignment error is controlled, leading to  
505 an  $O(1/\sqrt{T})$  convergence rate with improved regret bounds over standard SGD. The explicit dependencies on  
506 parameters (such as  $\epsilon$ ,  $\rho$ , and  $\varepsilon$ ) bridge the gap between theory and empirical performance.

## 507 B.7 Discussion and Practical Relevance

508 The final convergence bound explicitly depends on parameters such as  $\varepsilon$  (which captures the alignment error),  $\rho$   
509 (the variance contraction factor), and  $\eta$ . Explicitly, we have

$$\frac{1}{T} \sum_{t=1}^T \|\nabla f(w_t)\|^2 \leq \frac{f(w_1) - f(w^*)}{\eta(1-\varepsilon)T} + \frac{L\eta}{2(1-\varepsilon)} \rho n. \quad (\text{B.42})$$

510 Our concentration lemmas (e.g., Lemma B.1) guarantee that for sufficiently large group sizes, the empirical  
511 estimates  $\mu_j$  and  $\sigma_j^2$  deviate from the true moments by at most an error that decays exponentially in  $|G_j|$ . This  
512 justifies the approximate unbiasedness assumption and confirms that, under practical conditions (as seen in our  
513 empirical experiments), the constants  $\varepsilon$  and  $\rho$  remain small. Similarly, by ensuring that  $\epsilon$  is chosen relative to  
514 the minimum variance  $\sigma_{\min}$  via Lemma B.3, even edge cases where  $\sigma_j$  is very small are handled gracefully,  
515 ensuring stable normalization.

## 516 B.8 Discussion of Key Technical Aspects

### 517 B.8.1 Unbiasedness of the Normalized Gradient

518 Since the normalization is a nonlinear transformation, in general we cannot expect

$$\mathbb{E}[\tilde{g}_t] = \nabla f(w_t) \quad (\text{B.43})$$

519 to hold exactly. In our analysis, we assume approximate unbiasedness is achieved when:

- 520 1. The group sizes  $|G_j|$  are large so that the empirical estimates  $\mu_j$  and  $\sigma_j^2$  converge to the true moments.
- 521 2. The gradient components within each group come from a symmetric (or nearly symmetric) distribution.

522 Under these conditions, any bias induced by the normalization is of order  $O(\delta)$  (with  $\delta$  small). Future work will  
523 provide more explicit concentration bounds on this bias.

524 **B.8.2 Justification of the Variance Contraction Factor**

525 We require that

$$\rho_j = \frac{\sigma_j^2}{(\sigma_j + \epsilon)^2} \leq \rho < 1. \quad (\text{B.44})$$

526 A sufficient condition is that, for some constant  $c > 0$ , we have

$$\epsilon \geq c \sigma_j \quad \text{for all } j. \quad (\text{B.45})$$

527 Then,

$$\sigma_j + \epsilon \geq (1 + c)\sigma_j \Rightarrow \frac{\sigma_j}{\sigma_j + \epsilon} \leq \frac{1}{1 + c}, \quad (\text{B.46})$$

528 and thus

$$\rho_j \leq \left( \frac{1}{1 + c} \right)^2 \equiv \rho. \quad (\text{B.47})$$

529 We include this formally as Lemma B.4

530 **B.8.3 Alignment Between Normalized and True Gradients**

531 A key step in our convex analysis is the approximation

$$\langle \nabla f(w_t), \tilde{g}_t \rangle \approx \|\nabla f(w_t)\|_2^2. \quad (\text{B.48})$$

532 Define the error vector  $\delta_t$  by

$$\tilde{g}_t = \nabla f(w_t) + \delta_t. \quad (\text{B.49})$$

533 Then,

$$\langle \nabla f(w_t), \tilde{g}_t \rangle = \|\nabla f(w_t)\|_2^2 + \langle \nabla f(w_t), \delta_t \rangle. \quad (\text{B.50})$$

534 Under the assumption that the normalized gradient is nearly aligned with  $\nabla f(w_t)$  (which holds if the deviations  
535 within each group are small relative to  $\mu_j$ ), we can show that

$$|\langle \nabla f(w_t), \delta_t \rangle| \leq \varepsilon \|\nabla f(w_t)\|_2^2, \quad (\text{B.51})$$

536 for some  $\varepsilon = O(\delta)$  with  $\delta \ll 1$ . A formal lemma with full details is deferred to future extensions.

537 **B.8.4 Handling Small Denominators and the Role of  $\epsilon$**

538 The normalization divides by  $\sigma_j + \epsilon$ . In cases where  $\sigma_j$  is very small, the update may be dominated by  $\epsilon$ . To  
539 prevent unwanted distortion, one can:

- 540 1. Assume a lower bound  $\sigma_j \geq \sigma_{\min} > 0$  or,
- 541 2. Choose  $\epsilon$  sufficiently small relative to the typical scale of  $\sigma_j$ ,
- 542 3. Or analyze the two regimes separately.

543 We provide a preliminary lemma showing that, if  $\sigma_j$  is extremely small but the deviations  $|[g_t]_i - \mu_j|$  are also  
544 small, then the normalized gradient remains bounded.

545 **Lemma B.5** (Robustness for Small  $\sigma_j$ ). *Assume that for each group  $G_j$ , either*

- 546 1.  $\sigma_j \geq \sigma_{\min} > 0$ , or
- 547 2. When  $\sigma_j < \sigma_{\min}$ , the deviations satisfy  $|[g_t]_i - \mu_j| \leq \beta \sigma_{\min}$  for some  $\beta > 0$ .

548 Then, by choosing  $\epsilon \leq \gamma \sigma_{\min}$  for a sufficiently small  $\gamma > 0$ , we have

$$|[g_t]_i| \leq \frac{\beta}{\gamma}, \quad \text{for } i \in G_j, \quad (\text{B.52})$$

549 ensuring that the update remains bounded.

550 **C Convergence Proof for Non-Convex Settings**

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

554 *Proof Sketch of Theorem C* Under Assumptions C.1 C.2 and the variance contraction/moment conditions  
555 (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

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

556

□

## 557 C.1 Assumptions and Preliminaries

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

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

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

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

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

563 with  $G > 0$ .

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

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

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

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

567 and the normalized gradient is defined as

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

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

570 **C.2 Proof of Convergence**

571 The Milo update is given by

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

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

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

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

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

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

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

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

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

577 for some  $\varepsilon$  (dependent on  $\epsilon/\sigma_{\min}$ ).

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

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

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

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

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

582 Rearrange Equation C.12 to obtain an inequality on the gradient norm:

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

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

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

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

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

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

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

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

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

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

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

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

591 **C.3 Summary of the Proof**

592 1. **Smoothness:** Using  $L$ -smoothness we derived

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

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

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

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

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

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

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

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

597 **C.4 Discussion**

598 The analysis shows that even in the non-convex setting, the normalized gradients provided by Milo yield a  
599 controlled descent of the objective  $F$ . The key technical tools are:

- 600 • **Smoothness:** Ensuring that the objective does not change too abruptly.
- 601 • **Variance Contraction:** The grouping and normalization (via Assumption B.3) reduce the effective  
602 noise in the gradient, leading to a tighter bound on  $\|\tilde{g}_t\|$ .
- 603 • **Alignment Control:** Lemma B.2 guarantees that the error introduced by normalization,  $\delta_t$ , is suffi-  
604 ciently small relative to the true gradient.

605 Together, these factors allow us to conclude that Milo reaches an approximate stationary point (where  $\|\nabla F(w_t)\|$   
606 is small) at a rate of  $O(1/\sqrt{T})$ , even when the objective function is non-convex.

607 **C.5 Conclusion**

608 We have provided a detailed, self-contained convergence proof for Milo in the non-convex setting. By leveraging  
609 the same normalization strategy that yields variance contraction in the convex case and carefully controlling the  
610 alignment error, we show that Milo achieves an  $O(1/\sqrt{T})$  convergence rate in terms of the squared gradient  
611 norm. This result provides theoretical support for the empirical success of Milo in both convex and non-convex  
612 scenarios.

613 **D Experimentation**

614 **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

615 **D.2 Libraries Used**

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

- 617 • **PyTorch** (v2.6.0; BSD 3-clause License; <https://pytorch.org>) [Paszke et al., 2019]
- 618 • **torchvision** (v0.21.0; BSD 3-clause License; <https://github.com/pytorch/vision>) [PyTorch  
619 Contributors, 2023]
- 620 • **seaborn** (v0.13.2; BSD License; <https://seaborn.pydata.org>) [Waskom, 2021]
- 621 • **matplotlib** (v3.10.0; PSF License; <https://matplotlib.org>) [Hunter, 2007]
- 622 • **NumPy** (v2.1.3; BSD License; <https://numpy.org>) [van der Walt et al., 2011]
- 623 • **NovoGrad** optimizer (preprint; Apache-2.0; <https://arxiv.org/abs/1905.11286>) [Ginsburg  
624 et al., 2019]
- 625 • **pandas** (v2.2.3; BSD License; <https://pandas.pydata.org>) [McKinney, 2010]
- 626 • **JSON module** (Python 3.10 stdlib; PSF License; <https://docs.python.org/3/library/json.html>) [Crockford, 2006]
- 627 • **Optuna** (v4.3.0; MIT License; <https://optuna.org>) [Akiba et al., 2019]
- 628 • **scikit-learn** (v1.6.1; BSD License; <https://scikit-learn.org>) [Pedregosa et al., 2011]
- 629 • **SciPy** (v1.15.1; BSD License; <https://scipy.org>) [Virtanen et al., 2020]
- 630 • **itertools** (Python 3.10 stdlib; PSF License; <https://docs.python.org/3/library/itertools.html>) [Python Software Foundation, 2023]
- 631 • **Gymnasium** (v1.1.0; MIT License; <https://gymnasium.farama.org>) [Brockman et al., 2016]
- 632 • **Transformers** (v4.30.0; Apache-2.0; <https://github.com/huggingface/transformers>) [Wolf  
633 et al., 2020]
- 634 • **train-lm-from-scratch** (MIT License; <https://github.com/FareedKhan-dev/train-lm-from-scratch>) [FareedKhan-dev, 2025]
- 635 • **NovoGrad-pytorch** (MIT License; <https://github.com/lonePatient/NovoGrad-pytorch>)  
636 [lonePatient, 2019]
- 637
- 638
- 639

640 **D.3 Datasets**

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

- 642 • **MNIST** (Public Domain; <http://yann.lecun.com/exdb/mnist/>) [LeCun et al., 1998]
- 643 • **CIFAR-10** (2009 Technical Report; MIT License; <https://www.cs.toronto.edu/~kriz/cifar.html>) [Krizhevsky and Hinton, 2009]
- 644 • **CIFAR-100** (2009 Technical Report; MIT License; <https://www.cs.toronto.edu/~kriz/cifar.html>) [Krizhevsky and Hinton, 2009]
- 645 • **The Pile** (Apache 2.0 License; <https://pile.eleuther.ai/>) [Gao et al., 2020]
- 646
- 647

648 **D.4 Error Bar Computation**

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

651 **1. Calculation**

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

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

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

- 654 • **95% Confidence Interval (CI):**

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

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

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

658 **2. Sources of Variability**

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

663 **3. Assumptions**

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

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

670 **D.5 Experiment 1: Loss Map Predictions**

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

675 **D.5.1 Himmelblau Function and Its Gradient**

676 The Himmelblau function is defined as:

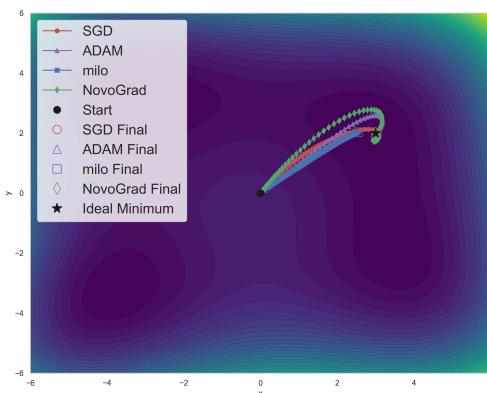
$$f(x_1, x_2) = (x_1^2 + x_2 - 11)^2 + (x_1 + x_2^2 - 7)^2. \quad (\text{D.1})$$

677 For the analytic gradient, we define:

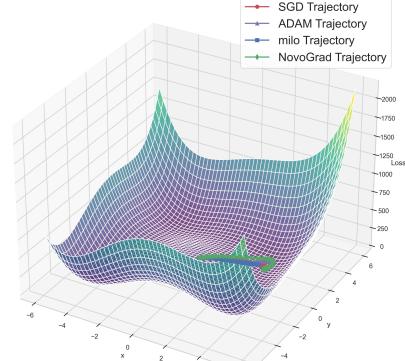
$$u = x_1^2 + x_2 - 11, \quad v = x_1 + x_2^2 - 7. \quad (\text{D.2})$$

678 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

679 **Experimental Setup**

- 680 • **Domain:**  $[-6, 6] \times [-6, 6]$   
 681 • **Initial Point:**  $(0, 0)$   
 682 • **Iterations:** 150

683 **Observations and Analysis**

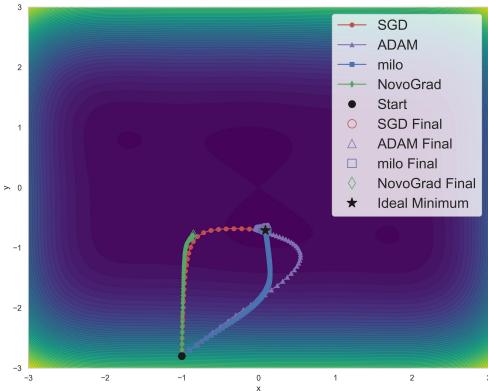
- 684 • **SGD:** The SGD trajectory follows a generally follows a more direct path, quickly descending into a  
 685 low-loss valley, curving towards the end, and stabilizing near the minimum.  
 686 • **Adam:** Adam has a less direct path, overshooting and eventually curving around to reach the objective.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

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

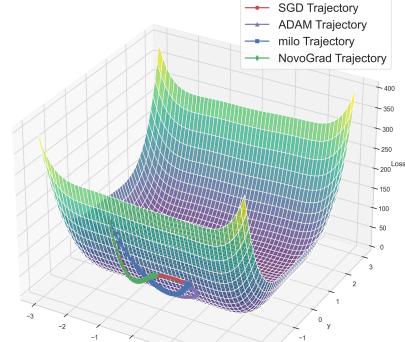
693 **D.5.2 Camel6 Function**

694 The Camel6 function is defined by:

$$f(x_1, x_2) = \left(4 - 2.1x_1^2 + \frac{x_1^4}{3}\right)x_1^2 + x_1x_2 + \left(-4 + 4x_2^2\right)x_2^2. \quad (\text{D.4})$$



(a) Camel6 Loss 2D



(b) Camel6 Loss 3D

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

695 **Experimental Setup**

- 696 • **Domain:**  $[-3, 3] \times [-3, 3]$   
 697 • **Initial Point:**  $(-1, -2.8)$   
 698 • **Iterations:** 150

699 **Observations and Analysis**

- 700 • **SGD:** The trajectory takes consistent steps, taking a largely indirect path to the local minima and  
 701 eventually converging to the global.  
 702 • **Adam:** Similar to SGD Adam takes an indirect path to the local minima, then shifts to a smaller step  
 703 size with a more direct path to the eventual global minima.  
 704 • **NovoGrad:** Novograd follows a similar path as SGD, however it comes to an early stop after it  
 705 reaches the local minima, never reaching the global.

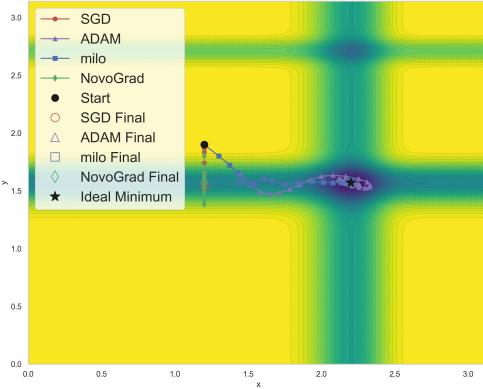
- 706 • **Milo:** With normalized updates, Milo takes consistent steps, initially following the same path as Adam,  
 707 but deviating towards the global minima for a more direct convergence than the other optimizers.

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

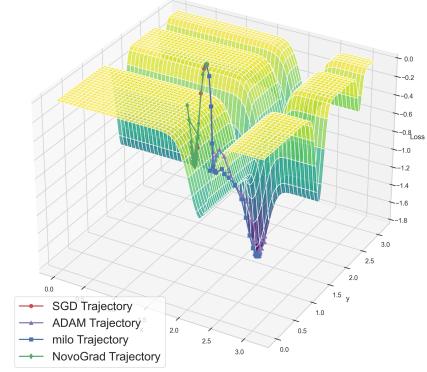
### 711 D.5.3 Michalewicz Function

712 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

### 713 Experimental Setup

- 714 • **Domain:**  $[-10, 10] \times [-10, 10]$
- 715 • **Initial Point:**  $(2.4, -4.3)$
- 716 • **Iterations:** 150

### 717 Observations and Analysis

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

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

### 729 D.5.4 De Jong Function 5

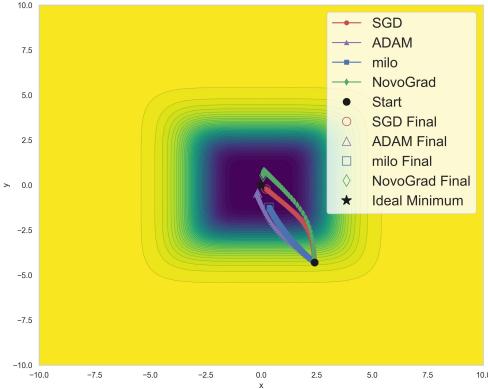
730 De Jong's fifth function is a multimodal test function defined for a 2-dimensional input vector  $x = (x_1, x_2)$ . It  
 731 is constructed using a fixed grid of constants and is known for its complex landscape with many local minima.  
 732 The function is given by:

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

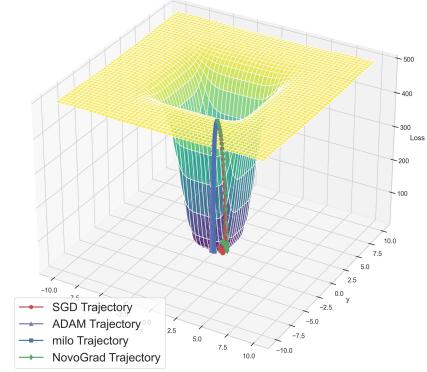
733 where the constants  $a_{1i}$  and  $a_{2i}$  form a  $5 \times 5$  grid of values from the set  $\{-32, -16, 0, 16, 32\}$ , specifically:

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

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



(a) DeJong Loss 2D



(b) DeJong Loss 3D

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

## 736 Experimental Setup

- 737 • **Domain:**  $[0, \pi] \times [0, \pi]$
- 738 • **Initial Point:**  $(1.2, 1.9)$
- 739 • **Iterations:** 250

## 740 Observations and Analysis

- 741 • **SGD:** Due to the steep descent structure of the function, SGD curves around the edge until it corrects  
742 towards the global minima, reaching it.
- 743 • **Adam:** Adam's adaptive updates help it follow a smoother descent, reaching the final global minima  
744 smoothly.
- 745 • **NovoGrad:** NovoGrad follows a similar path around the edge as SGD, eventually overshooting  
746 slightly before reaching the global minima.
- 747 • **Milo:** Milo's normalized update rule maintains a consistent step size even when the gradients are  
748 small, allowing steady progress; however, it stops right as the bottom of the basin is reached, not  
749 continuing further as the other optimizers do.

## 750 Key Takeaway:

### 751 D.5.5 Experiment 1 Analysis

#### 752 Optimizer Comparison

- 753 • **Early Stage Behavior:**
  - 754 – **SGD:** May take large initial steps but risks overshooting and stalling on plateaus.

- 755           – **Adam:** Adapts the step size effectively, resulting in a rapid descent and smoother convergence,  
 756           but also commonly overshoots.
- 757           – **NovoGrade:** NovoGrad has similar behavior as Adam and SGD, however it displays issues  
 758           with reaching the global minima when there are many local minima as well.
- 759           – **Milo:** Uses normalized updates to maintain consistent step sizes, ensuring steady progress even  
 760           when gradients are small.
- 761           • **Navigation of Plateaus and Local Minima:**
- 762           – SGD may get trapped in local minima without additional momentum or careful scheduling.
- 763           – Adam can sometimes reduce its step size too quickly, causing temporary plateaus.
- 764           – NovoGrad performs well when in environments with few local minima, but shows issues with  
 765           handling higher dimensionality functions.
- 766           – Milo excels at escaping regions with vanishing gradients by enforcing a constant update magni-  
 767           tude.
- 768           • **Final Convergence and Stability:**
- 769           – Adam and Milo tend to converge to lower loss values than SGD and NovoGrad.
- 770           – Overall, Milo demonstrates both robust exploration in complex landscapes and stable conver-  
 771           gence.

## 772      Practical Takeaways

- 773           • **Function-Specific Behavior:**
- 774           – In multimodal functions (Himmelblau, Michalewicz), normalization helps maintain progress.
- 775           – In functions with shallow basins (Michalewicz), Milo’s constant step size aids in escaping local  
 776           traps.
- 777           • **Optimizer Selection:**
- 778           – No single optimizer is universally best. Milo offers a balanced approach when gradient scales  
 779           vary widely and consistent progress is desired.

780      Overall the experiments reveal that while SGD, Adam, NovoGrad, and Milo each have their strengths, the  
 781      proposed Milo optimizer, with its normalized, group-based updates, can offer stable and consistent progress in  
 782      challenging optimization landscapes. In particular, Milo shows promise in escaping flat regions and avoiding  
 783      overshooting due to its consistent step size. However, as shown with DeJong, the consistent step size can  
 784      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 sam- pling 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 augmenta- tions)

785 **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 <sub>LW</sub>	– (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 <sub>LW</sub>	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

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

787 This experiment evaluates a convex multi-class logistic regression using the MNIST Dataset ([LeCun et al.,  
788 1998]), just as was done in the study that proposed Adam. The goal is to compare optimizers with respect to  
789 training loss, convergence speed, and test accuracy. The model uses L2 regularization with a single linear layer  
790 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 ( $28 \times 28$ )
Linear	Single fully connected layer mapping 784 to 10 (number of classes)

Table 8: Logistic Regression Training Results

<b>Optimizer</b>	<b>Epoch Group</b>	<b>Loss</b>	<b>Accuracy (%)</b>	<b>F1 Score</b>
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
$\text{Milo}_{\text{LW}}$	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)

<b>Optimizer</b>	<b>Avg Runtime (s)</b>	<b>Avg Memory (MB)</b>
MILLO	32.51	63.5
$\text{MILLO}_{\text{LW}}$	32.66	2.0
SGD	29.71	0.0
ADAGRAD	29.76	0.0
ADAMW	30.03	16.4
NOVOGRAD	32.16	4.1

### 791 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.
MILLO	$\text{MILLO}_{\text{LW}}$	0.350833	0.304008	$\text{MILLO}_{\text{LW}}$	8.13499e-08	***
MILLO	SGD	0.350833	0.306714	SGD	3.62048e-06	***
MILLO	ADAMW	0.350833	0.287661	ADAMW	9.31980e-07	***
MILLO	ADAGRAD	0.350833	0.481010	MILLO	9.90140e-11	***
MILLO	NOVOGRAD	0.350833	0.309194	NOVOGRAD	4.23780e-06	***
$\text{MILLO}_{\text{LW}}$	SGD	0.304008	0.306714	$\text{MILLO}_{\text{LW}}$	1.32372e-01	
$\text{MILLO}_{\text{LW}}$	ADAMW	0.304008	0.287661	ADAMW	9.06983e-05	***
$\text{MILLO}_{\text{LW}}$	ADAGRAD	0.304008	0.481010	$\text{MILLO}_{\text{LW}}$	2.49227e-13	***
$\text{MILLO}_{\text{LW}}$	NOVOGRAD	0.304008	0.309194	$\text{MILLO}_{\text{LW}}$	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 <sub>LW</sub>	0.912861	0.916955	MILO <sub>LW</sub>	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 <sub>LW</sub>	SGD	0.916955	0.914048	MILO <sub>LW</sub>	1.56224e-02	*
MILO <sub>LW</sub>	ADAGRAD	0.916955	0.876378	MILO <sub>LW</sub>	2.24596e-10	***
MILO <sub>LW</sub>	ADAMW	0.916955	0.919612	ADAMW	9.94675e-03	**
MILO <sub>LW</sub>	NOVOGRAD	0.916955	0.913075	MILO <sub>LW</sub>	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	***

## 792 D.8 Experiment 3: Multilayer Neural Networks

793 The objective of this experiment is to test the optimizer on a non-convex problem by using a multilayer perception.  
 794 For this study, our model choice was driven by the previous Adam publication, using two fully connected layers  
 795 with 1000 ReLU units each, with dropout and weight decay on a minibatch of 128, on the MNIST image dataset  
 796 ([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

<b>Optimizer</b>	<b>Epoch Group</b>	<b>Loss</b>	<b>Accuracy (%)</b>	<b>F1 Score</b>
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
$\text{Milo}_{\text{LW}}$	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)

<b>Optimizer</b>	<b>Avg Duration (s)</b>	<b>Avg Memory (MB)</b>
MILO	34.37	18.4
$\text{MILO}_{\text{LW}}$	36.05	0.0
SGD	30.28	0.0
ADAGRAD	30.74	0.0
ADAMW	30.87	0.0
NOVOGRAD	34.24	0.0

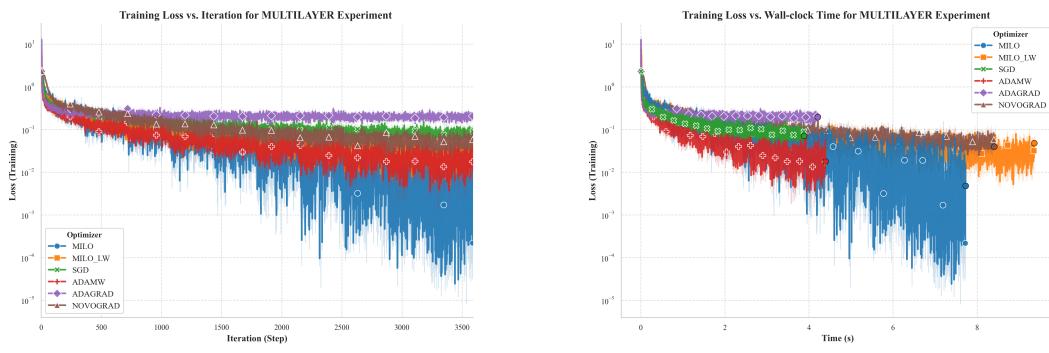


Figure 10: Comparison of MLP training cost by iterations and wall time.

797 **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 <sub>LW</sub>	0.225943	0.150747	MILO <sub>LW</sub>	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 <sub>LW</sub>	SGD	0.150747	0.104657	SGD	4.70080e-04	***
MILO <sub>LW</sub>	ADAMW	0.150747	0.121495	ADAMW	1.64005e-03	**
MILO <sub>LW</sub>	ADAGRAD	0.150747	0.215592	MILO <sub>LW</sub>	1.31936e-05	***
MILO <sub>LW</sub>	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 <sub>LW</sub>	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 <sub>LW</sub>	SGD	0.968024	0.971065	SGD	1.56224e-02	**
MILO <sub>LW</sub>	ADAMW	0.968024	0.972800	ADAMW	6.94322e-03	**
MILO <sub>LW</sub>	ADAGRAD	0.968024	0.941102	MILO <sub>LW</sub>	3.03333e-05	***
MILO <sub>LW</sub>	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	***

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

799 This experiment evaluates a deep Residual Network, specifically ResNet-18 [He et al., 2016], on the CIFAR-100  
800 dataset [Krizhevsky and Hinton, 2009] to test optimizer performance in a finer-grained classification setting.  
801 Compared to CIFAR-10, CIFAR-100 has 100 target classes, which stresses both the representation capacity of  
802 the network and the optimizer's ability to generalize. We use the standard ResNet-18 architecture, modifying  
803 only the final linear layer to output 100 logits. Input images are augmented with 4-pixel zero-padding, then  
804 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

<b>Layer</b>	<b>Description</b>
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

<b>Optimizer</b>	<b>Epoch Group</b>	<b>Loss</b>	<b>Accuracy (%)</b>	<b>F1 Score</b>
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
$\text{Milo}_{\text{LW}}$	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)

<b>Optimizer</b>	<b>Avg Duration (s)</b>	<b>Avg Memory (MB)</b>
MILO	148.53	59.4
$\text{MILO}_{\text{LW}}$	182.82	2.0
SGD	130.56	14.3
ADAGRAD	130.69	0.0
ADAMW	131.10	14.3
NOVOGRAD	172.25	4.1

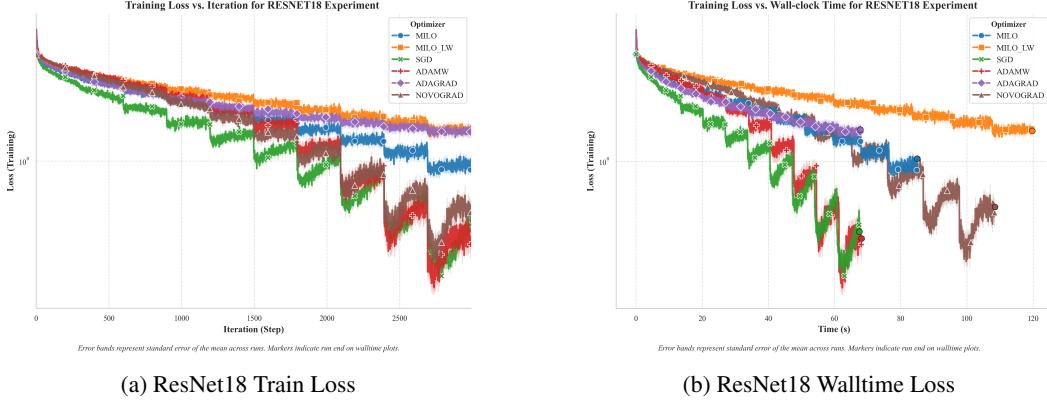


Figure 11: Comparison of ResNet18 training cost by iterations and wall time.

### 805 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 <sub>LW</sub>	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 <sub>LW</sub>	SGD	2.11451	2.16893	MILO <sub>LW</sub>	0.081725	
MILO <sub>LW</sub>	ADAGRAD	2.11451	2.48136	MILO <sub>LW</sub>	1.24390e-06	***
MILO <sub>LW</sub>	ADAMW	2.11451	3.38785	MILO <sub>LW</sub>	6.11940e-06	***
MILO <sub>LW</sub>	NOVOGRAD	2.11451	2.29087	MILO <sub>LW</sub>	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 <sub>LW</sub>	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 <sub>LW</sub>	SGD	0.441000	0.454009	SGD	0.050683	
MILO <sub>LW</sub>	ADAGRAD	0.441000	0.365256	MILO <sub>LW</sub>	1.69300e-06	***
MILO <sub>LW</sub>	ADAMW	0.441000	0.411299	MILO <sub>LW</sub>	0.024678	*
MILO <sub>LW</sub>	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	***

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

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

Table 22: RL Experiment Training Parameters

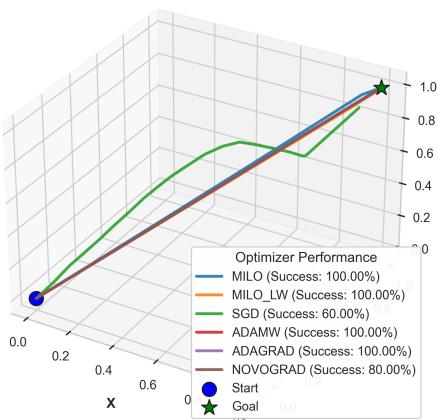
Parameter	Value
Learning Rate	0.0005
Gamma (Discount Factor)	0.99
Episodes	1000
Log Interval	250 episodes
Initial Exploration	0.2
Final Exploration	0.01
Exploration Decay	0.998
Gradient Clip	0.5
Distance Threshold	0.05
Random Goal	True
Max Steps per Episode	500
Start Position	[0.0, 0.0, 0.0]
Goal Position	[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"> <li>• Linear (<math>256 \rightarrow 128</math>)</li> <li>• LayerNorm + ReLU</li> <li>• Dropout (0.1)</li> <li>• Linear (<math>128 \rightarrow 256</math>)</li> <li>• LayerNorm</li> <li>• Skip connection</li> </ul>
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).

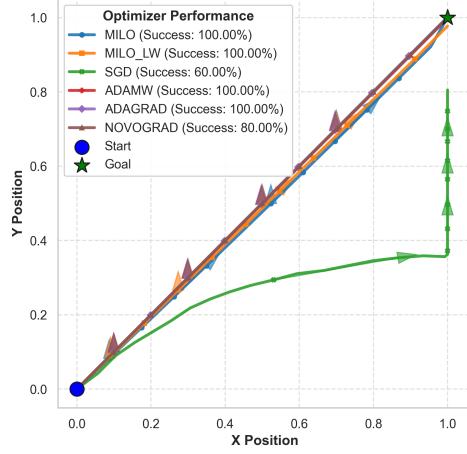
818 **D.10.1 Set Result:**

Average Agent Trajectories by Optimizer



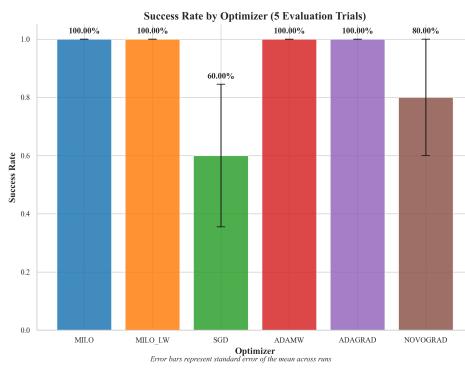
(a) 3D Agent Tractories

Top-Down View of Average Agent Trajectories



(b) 2D Agent Tractories

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



(a) RL Set Success Rates



(b) Reward Comparision

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

819 As shown in Figure 12 and Figure 13 each of the optimizers has a high success rate of reaching the set goal  
 820 with no random initialization or generation of the goal. The optimizers with the highest performance proved  
 821 to be AdamW, AdaGrad, Milo, and Milo<sub>LW</sub>, each achieving a 100% accuracy. Overall, the set experiment  
 822 demonstrates that each of the models can effectively learn the patterns generated in this 3D space.

823 **D.10.2 Randomized Result:**

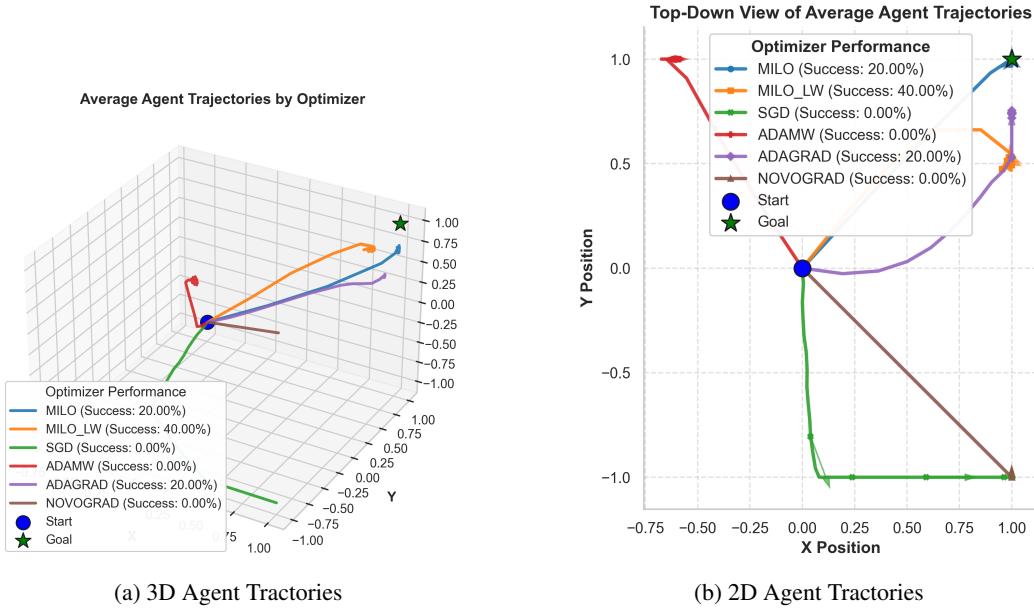


Figure 14: 3D and 2D visuals of set RL agent trajectories.

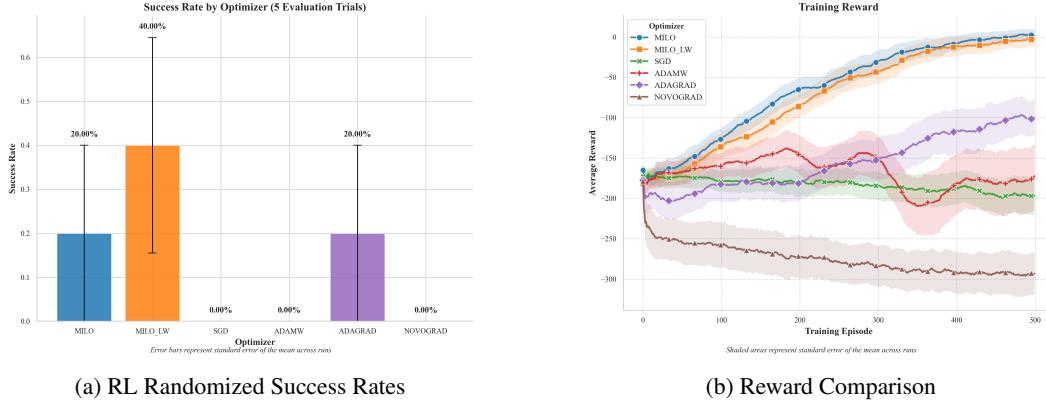


Figure 15: Comparison of randomized success rates by optimizer and average reward by episode.

824 The results of the randomly generated start and goal differ drastically from the set, with each optimizer performing  
 825 worse than before. However, each optimizer retains the relative performance shown previously regarding their  
 826 performance against one another. As shown in Figure 14a, Figure 14b, and Figure 15a, Milo<sub>LW</sub> has the highest  
 827 accuracy, with a final success rate of 40%, being a large improvement over the success rates of the other  
 828 optimizers. This performance is also shown in Figure 15b, with Milo<sub>LW</sub> demonstrating a higher reward than the  
 829 other optimizers.

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

831 LLMs have grown significantly in popularity in recent years, with many innovations made in model architecture  
 832 and optimization. This study explores the performance of our optimizer in a pretraining setting, using a small-  
 833 scale Transformer model trained from scratch on a processed subset of The Pile dataset [Gao et al., 2020]. The  
 834 model weights are randomly initialized at the start of training, and no pre-trained parameters are used at any  
 835 stage. This setup allows us to evaluate optimizer behavior in a controlled language modeling task, where both  
 836 model capacity and training steps are deliberately constrained to enable rapid comparison across runs.

837 To support this investigation, we adapt the open-source training framework of [FareedKhan-dev \(2025\)](#), which  
 838 implements a simplified version of the Transformer architecture introduced by [Vaswani et al. \(2017\)](#). This  
 839 architecture includes core components such as multi-head self-attention, layer normalization, and position-wise  
 840 feed-forward networks, implemented natively in PyTorch. We extend the original implementation to support  
 841 multiple optimizers and perform controlled training and evaluation across several runs. Visual results can be  
 842 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

843 The results in this experiment, as shown in Figure 4a and Figure 4b, demonstrate that Milo and Milo<sub>LW</sub>  
 844 demonstrate comparable performance to some of the industry standard used for training LLMs such as Adam  
 845 and AdamW, with Milo<sub>LW</sub> and Milo achieving the lowest recorded loss and fastest reduction of loss as well.  
 846 Additionally, since this was a more complex task, the wall time of each optimizer was logged as well, with Milo  
 847 being comparable to Adam and AdamW, whereas Milo<sub>LW</sub> and NovoGrad took a longer training time.

## 848 D.12 Experiment 7: Ablation Study

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

853 **Experimental Setup.** We evaluate two optimizer variants: standard Milo and its layer-wise variant Milo<sub>LW</sub>.  
 854 For each model–dataset pair, we establish a baseline configuration (see `config.py`) and systematically vary one  
 855 hyperparameter at a time. Parameters under test include:

- 856 • **Learning Rate (lr):** {0.0001, 0.001, 0.005, 0.01}
- 857 • **Momentum (momentum):** {0.0, 0.5, 0.9}
- 858 • **Weight Decay (weight\_decay):** {0, 1e-4, 1e-3}
- 859 • **Clipping Norm (clip\_norm):** {None, 1.0, 5.0}
- 860 • **Epsilon (eps):** {1e-8, 1e-5, 1e-2}
- 861 • **Scale-Aware Blending (scale\_aware, scale\_factor):** {True, False}, {0.1, 0.2, 0.5} (Milo<sub>LW</sub> only)

862 Each configuration is run for RUNS\_PER\_OPTIMIZER independent trials to ensure statistical robustness.  
 863 Results are recorded under structured directories and exported as CSV metrics and loss/accuracy curves.

864 **Metrics.** We measure:

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

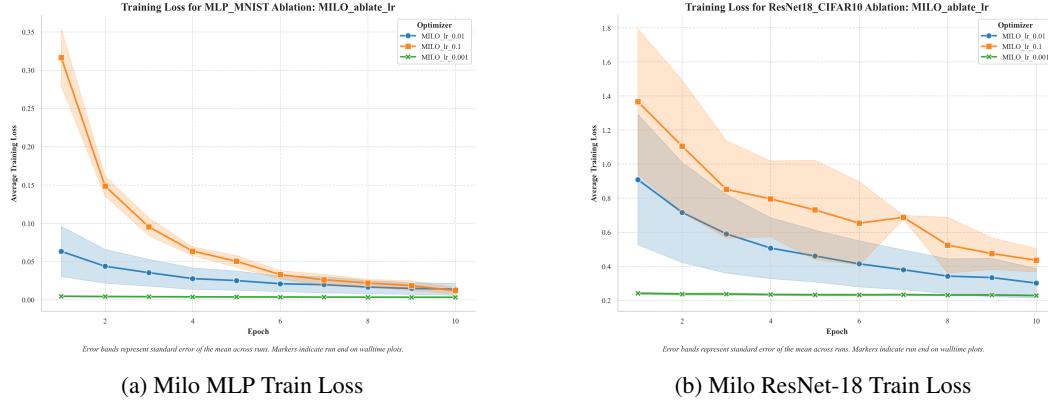


Figure 16: Comparison of Milo LR Ablation in MLP and ResNet-18.

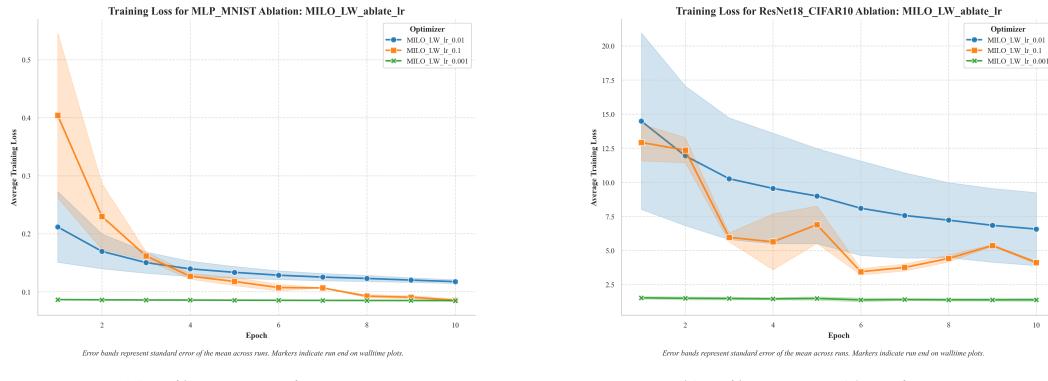


Figure 17: Comparison of Milo\_LW LR Ablation in MLP and ResNet-18.

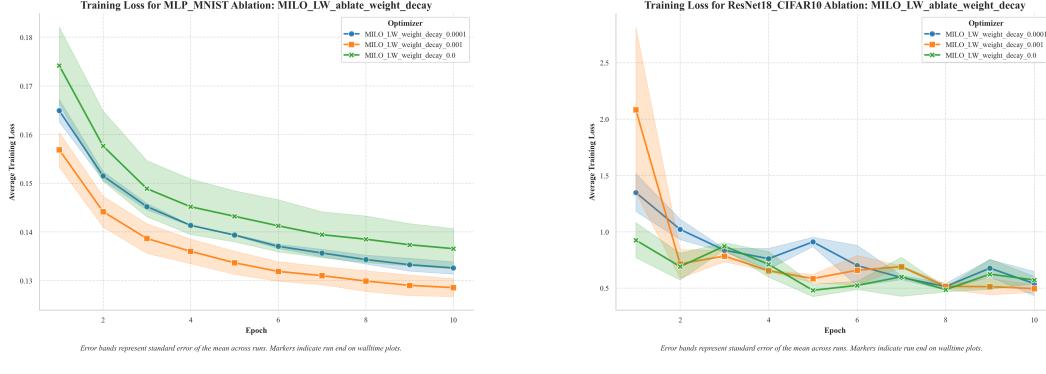


Figure 18: Comparison of Milo\_LW Weight Decay Ablation in MLP and ResNet-18.

### 868 Key Findings.

- 869 • **Learning Rate** remains the most critical factor: optimal values vary by task and model, with Milo and  
870 Milo<sub>LW</sub> exhibiting similar sensitivity profiles.
- 871 • **Weight Decay** and **Epsilon** show moderate, task-dependent effects; safe defaults (0 and  $10^{-5}$ ) yield  
872 robust performance.
- 873 • **Momentum** and **Clipping Norm** have marginal impact on final accuracy but can improve early  
874 stability.
- 875 • **Scale-Aware Blending**: Contrary to our hypothesis, blending (Milo<sub>LW</sub>) did not consistently outper-  
876 form standard Milo. In MLP/MNIST and DeepCNN/CIFAR-10, the best Milo<sub>LW</sub> (scale\_factor=0.5)  
877 still trailed standard Milo by 0.2–0.5% accuracy.

878 **Implications.** These results suggest that Milo’s group normalization and adaptive components are robust  
879 across a wide range of hyperparameters, reducing the tuning burden relative to conventional optimizers. However,  
880 the scale-aware blending mechanism may require further refinement or alternative blending strategies to realize  
881 its theoretical benefits in practice.