
Stepsize anything: A unified learning rate schedule for budgeted-iteration training

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 The expanding computational costs and limited resources underscore the critical
2 need for budgeted-iteration training, which aims to achieve optimal learning
3 within predetermined iteration budgets. While learning rate schedules fundamen-
4 tally govern the performance of different networks and tasks, particularly in
5 budgeted-iteration scenarios, their design remains largely heuristic, lacking the-
6 oretical foundations. In addition, the optimal learning rate schedule requires exten-
7 sive trial-and-error selection, making the training process inefficient. In this work,
8 we propose the Unified Budget-Aware (UBA) schedule, a theoretically grounded
9 learning rate schedule that consistently outperforms commonly-used schedules
10 among diverse architectures and tasks under different constrained training budgets.
11 First, we bridge the gap by constructing a novel training budget-aware optimiza-
12 tion framework, which explicitly accounts for the robustness to landscape curva-
13 ture variations. From this framework, we derive the UBA schedule, controlled
14 by a single hyper-parameter φ that provides a trade-off between flexibility and
15 simplicity, eliminating the need for per-network numerical optimization. More-
16 over, we establish a theoretical connection between φ and the condition number,
17 adding interpretation and justification to our approach. Besides, we prove the con-
18 vergence for different values of φ . We offer practical guidelines for its selection
19 via theoretical analysis and empirical results. Extensive experimental results show
20 that UBA *consistently surpasses* the commonly-used schedules across diverse vi-
21 sion and language tasks, spanning network architectures (e.g., ResNet, OLMo)
22 and scales, under different training-iteration budgets.

1 Introduction

23 Deep learning has achieved remarkable success across various across a wide range of domains, in-
24 cluding computer vision and natural language processing. However, despite continual advancements
25 in hardware technologies [44, 52], the training cost of neural networks has increased dramatically
26 due to the growing scale of models and datasets [4, 8, 46, 47]. As a result, resource constraints,
27 including computational power, memory, energy consumption, and time budgets, are emerging as
28 significant bottlenecks in the training process [40, 57]. These challenges highlight the pressing need
29 for budgeted training, which aims to achieve optimal model performance under fixed hardware and
30 limited time.

32 While existing budgeted training studies broadly address resource efficiency, a critical yet under-
33 explored direction within budgeted training is achieving the best possible model performance under
34 strictly fixed iteration constraints. This scenario is common and practically significant where practi-
35 tioners work under limited computational or time budgets [29], and in extreme cases, models have to
36 be completed within a few training iterations due to resource exhaustion. To formalize this specific

37 research problem, we introduce the term ‘budgeted-iteration training’ , distinguishing it from the
38 broader scope of budgeted training.

39 Budgeted-iteration training has received growing attention in research, given its significant real-
40 world applicability. Several studies have developed relevant techniques that align with its goals.
41 Smith et al. [42] propose the cyclical learning rate schedule (CLR), improving accuracy in fewer
42 iterations without tuning [43]. Li et al. [29] introduce budget-aware adaptations for existing learning
43 rate schedules. Chen et al. [5] propose a novel learning rate schedule called Reflected Exponential
44 (REX). These approaches are primarily based on learning rate designing. Learning rate scheduling
45 highlights key advantages: (i) It plays a critical role in general training of diverse neural architectures
46 across tasks. (ii) It has demonstrated suitable and competitive in fixed training iteration budgets [19].
47 (iii) It is plug-and-play, requiring minimal adjustments to the underlying model architecture, which
48 makes them easily adaptable to various deep learning frameworks. Leveraging these advantages, we
49 adopt the learning rate design approach for budgeted-iteration training.

50 Despite their advantages, most learning rate schedules, whether tailored for budgeted-iteration
51 training or standard training performance, are still heuristic and lack rigorous theoretical grounding.
52 In addition, existing schedules typically rely on manually designed rules or empirical tuning.
53 Consequently, selecting an optimal schedule often involves extensive trial-and-error, incurring
54 substantial cost in terms of time and computation [33].

55

56 *A natural question arises: Does there exist a theoretically grounded, unified schedule that*
57 *eliminates heuristic selection while maintaining robust performance across tasks, networks, scales*
58 *and training budgets?*

59 In this paper, we provide an affirmative answer by proposing a theoretically grounded schedule. The
60 proposed learning rate schedule should consistently outperforms existing schedules among diverse
61 architectures and tasks under different constrained training budgets. By doing so, it avoids choosing
62 suitable learning rate schedule after multiple trials for network training.

63 To achieve this, we first bridge the gap by constructing a unified budget-aware training optimization
64 framework, which incorporates the robustness to landscape curvature variations induced by data dis-
65 tribution, sampling, network architectures and optimization. Then we obtain numerical solutions by
66 gradient projection methods. To eliminate the need for repeated numerical optimization when apply-
67 ing our method to different networks, we propose a universal parametric function that approximates
68 numerical solutions. We nominate the resulting schedule Unified Budget-Aware (UBA) schedule. It
69 requires tuning only a single hyper-parameter φ , reducing the overhead of per-network numerical
70 optimization. Moreover, we establish a theoretical connection between φ and the condition number,
71 adding interpretation and optimization difficulty-aware theoretical grounding. Besides, we prove the
72 convergence for different values of φ . These theoretical analysis along with empirical results offer
73 practical guidelines for φ selection. We evaluate UBA through comprehensive experiments across
74 vision and language tasks, spanning diverse architectures and iteration budgets. Specifically, for vi-
75 sion tasks, the UBA schedule demonstrates *consistent superiority* over baselines across all evaluated
76 datasets and model scales under varying training iterations. For language tasks, we validate the ef-
77 fectiveness of UBA through extensive benchmarks with OLMo model(36M, 73M, 150M and 300M
78 parameters). Results show that UBA achieves state-of-the-art performance across on approximately
79 half of the benchmarks, and consistently outperforms baselines on the average scores.

80 **Main contributions:**

- 81 1. We construct a unified budget-aware training optimization framework that inherently adapts
82 to landscape curvature variations, enhancing training robustness of leaning rates.
- 83 2. We propose the Unified Budget-Aware (UBA) schedule from our constructed optimization
84 problem, controlled by a single hyper-parameter φ that provides a trade-off between flexi-
85 bility and simplicity, i.e. adaptive curvature adjustment and minimal tuning cost.
- 86 3. We prove the convergence under φ and derive practical guidelines for its selection through
87 analysis and experiments. Besides, theoretical analysis and empirical results show that φ
88 is related to optimization difficulty such as condition number.

89 4. We perform experiments and demonstrate that UBA surpasses the commonly-used schedules
90 across diverse vision and language tasks, spanning network architectures (e.g., ResNet,
91 OLMo) and scales, under different training-iteration budgets.

92 2 Related work

93 **Budgeted training** Researchers face significant challenges in achieving optimal model performance
94 under fixed hardware and limited time. To address these challenges, the concept of budgeted
95 training has gained increasing attention, exploring techniques including computation efficiency,
96 model compression, training stability and convergence improvement [40]. It focuses on: (i) em-
97 phasizing the allocation of resources, such as the balance between the model size and the amount
98 of data [2, 6, 23, 27]. (ii) finding optimal configurations or improving performance within the given
99 compute or time budget, such as memory efficiency and computation reduction [16, 25, 30, 37], op-
100 timization learning rate schedules [5, 29, 42, 43], batch size [16] and other weight averaging method
101 [25, 26].

102 Within the context of budgeted training, a key area that remains under-explored is achieving the best
103 possible model performance within a fixed number of iterations, i.e. ‘budgeted-iteration training’. In
104 this domain, learning rate scheduling based methods are particularly aligned with the objectives of
105 budgeted-iteration training. Smith et al. [42] propose a new learning rate schedule, named cyclical
106 learning rates (CLR). It improves accuracy in fewer iterations without tuning and is relevant to super-
107 convergence phenomenon [43]. Li et al. [29] introduce an alternative setting of existing learning
108 rate schedules for budgeted training. Chen et al. [5] propose the reflected exponential schedule
109 (REX) via a profile and sampling fashion. Learning rate based approaches achieve robust and high-
110 performing results under various lengths of training iterations, which corresponds to our purpose in
111 this work, i.e. budgeted-iteration training [33]. Besides, this approach is plug-and-play, requiring no
112 substantial alterations to the underlying model structure, making it readily adaptable to various deep
113 network frameworks. Therefore, we explore the proper learning rate schedule to achieve budgeted-
114 iteration training.

115 **Learning rate schedule** The learning rate plays a pivotal role in controlling the optimization pro-
116 cess during network training. The common scheme is the step decay schedule. A typical instance
117 decreases the learning rate by a decaying scalar 0.1 after 50% epochs and by a decaying scalar 0.01
118 after 75% epochs [20]. Then Loshchilov et al. [32] observe that sharp decreases may prevent models
119 from escaping local minima and propose the cosine schedule function, which is the most popular
120 schedule for language model pretraining. Although some schedules include the CLR [42], REX [5],
121 Warmup-Stable-Decay (WSD) [24] and schedule from multi-power law [33], there is no consensus
122 on the optimal choice. In addition, the detail can be found in Appendix F, including some works
123 focus on adaptive learning rate methods.

124 Learning rate design is not only significant in general training, but also critical in budgeted-iteration
125 training which still remains a topic of debate. Some analyses advocate for small, constant learning
126 rates to ensure stability and convergence [11]. On the contrast, one prevailing hypothesis suggests
127 that large learning rates may facilitate crossing over sharp local minima in the optimization land-
128 scape [55]. Despite the lack of comprehensive theoretical explanations, a range of learning rate
129 schedules inspired by the above analyses as heuristic guidelines has been widely adopted in practice,
130 using variable learning rates to budgeted-iteration training [5, 29]. In this work, we explore learning
131 rate schedule from optimization problem tailored to budgeted-iteration training, aiming to balance
132 iteration budget constraints and generalization.

133 3 Budgeted-iteration training

134 3.1 Finite optimization under limited training iterations

135 To design a one-size-fits-all learning rate schedule, we construct a robust optimization model of
136 learning rates across varying training conditions (see more conceptual illustration in Appendix A).
137 Specifically, we aim to guarantees minimal loss within constrained training iterations under the
138 worst-case conditions, proposing a budget-iteration-aware framework for learning rate optimization.

139 **Definition (Finite optimization):** Let $f(W, D)$ denote the function parameterized by a given neural
 140 network with parameters $W \in \mathbb{R}^N$ on the dataset D , ξ denotes the data sampling on the dataset
 141 D , and let \mathcal{F} be a function class. Let \mathcal{L} be the loss function. Let η_t be the learning rate at the t -th
 142 iteration, T be a maximum number of training iterations, and t be the current learning step. A finite
 143 optimization is

$$\begin{aligned} & \min_{\eta_1, \eta_2, \dots, \eta_{T-1}} \max_{f \in \mathcal{F}} \mathcal{L}(f(W_T, \xi)) \\ & \text{s.t. } W_{t+1} = W_t - \eta_t \nabla \mathcal{L}(f(W_t, \xi)) \\ & \quad t = 0, 1, 2, \dots, T-1 \end{aligned} \tag{1}$$

144 In the optimization model 1, the constraint represents the stochastic gradient descent process. The
 145 maximizing of $\mathcal{L}(f(W_T, \xi))$ represents the worst-case among the training process on the net-
 146 work f . Then it minimizes the worst-case loss within given iterations, embodying its budget-
 147 aware property. By formulating the problem as a min-max optimization, we identify learning rates
 148 $\eta_t (t = 1, 2, \dots, T-1)$ that are resilient to the uncertainties introduced by different training config-
 149 uration, uniformly throughout the optimization trajectory.

150 The challenge lies in characterizing the f within the optimization process. f is primarily deter-
 151 mined by variations in parameter configuration, datasets characteristics, batch ordering and network
 152 architecture. From an optimization standpoint, we assume that the characteristics of f shaped by
 153 these factors can be captured by the loss landscape of f . Therefore, we can analyze f by ap-
 154 proximating its loss using a quadratic expansion around nearby strict local optima. Specifically,
 155 during the optimization process, the loss surface in the vicinity of the optimization trajectory can
 156 be approximated by sequence of strict local optima (or at least by one optimum), denoted as
 157 $\bar{W}^{(k)} \quad (k = 1, 2, \dots, K) \quad K \in \mathbb{Z}^+$, with the final optimum represented as \bar{W}^K . This approach
 158 enables us to capture the key features of the loss surface near these points. Therefore, the trajec-
 159 tory of optimization is impacted by the characteristics of the nearby strict local optima, since the
 160 optimization process is inherently shaped by the loss landscape and the key features of the surface
 161 are captured by these optima. Consequently, we can derive the learning rate within the optimization
 162 model by the information of these nearby strict local optima.

163 According to the second-order necessary condition for the strict minimum $\bar{W}^{(k)}$, this approximation
 164 is achieved through the positive semi-definite Hessian matrix $H_f^{(k)}(\xi) \in \mathbb{R}^{N \times N}$. In addition, the
 165 optimization problem (1) will be programmed sequentially. Then objective function of the optimiza-
 166 tion problem (1) can be reformulated as $\mathcal{L}(f(\bar{W}^{(k)}, \xi)) + \frac{1}{2}(W_{T_{k+1}} - \bar{W}^{(k)})^\top H_f^{(k)}(\xi)(W_{T_{k+1}} -$
 167 $\bar{W}^{(k)}) \quad (k = 1, 2, \dots, K)$. Given that the networks under consideration possess sufficient capac-
 168 ity to fit the data, the loss at the optimal point for different f can be made small enough. Thus the
 169 term $\mathcal{L}(f(\bar{W}^{(k)}, \xi))$ can be reasonably neglected in our analysis. We obtain the following sequential
 170 optimization problem.

$$\begin{aligned} & \min_{\eta_{1+T_k}, \eta_{2+T_k}, \dots, \eta_{T_{k+1}}} \max_{f(\xi)} \frac{1}{2}(W_{T_{k+1}} - \bar{W}^{(k)})^\top H_f^{(k)}(\xi)(W_{T_{k+1}} - \bar{W}^{(k)}) \\ & \text{s.t. } W_{t+1} = W_t - \eta_t H_f^{(k)}(\xi)(W_t - \bar{W}^{(k)}) \\ & \quad t = T_k + 1, T_k + 2, \dots, T_{k+1} \quad (k = 1, 2, \dots, K-1) \end{aligned} \tag{2}$$

171 where $T_1 = 0$ and $T_K = T-1$.

172 By the first constraint $W_{t+1} = W_t - \eta_t H_f^{(k)}(\xi)(W_t - \bar{W}^{(k)})$, we can obtain (see Appendix A for
 173 derivation)

$$\left\| W_{T_{k+1}} - \bar{W}^{(k)} \right\|_2^2 \leq \max_{\lambda_l^{(k)} \leq \lambda_i^{(k)} \leq \lambda_u^{(k)}} \left[\prod_{t=1+T_k}^{T_{k+1}} (1 - \eta_t \lambda_i^{(k)}) \right]^2 \left\| W_{T_k} - \bar{W}^{(k)} \right\|_2^2 \tag{3}$$

174 where $\lambda_i^{(k)} \quad (i = 1, 2, \dots, N)$ are the eigenvalues of $H_f^{(k)}(\xi)$ around the k -th optimum, which
 175 satisfy $0 < \lambda_l^{(k)} \leq \lambda_i^{(k)}(f(\xi)) \leq \lambda_u^{(k)}$ for all i , $u_i^{(k)} \quad (i = 1, 2, \dots, N)$ denote N linearly in-
 176 dependent eigenvectors and $s_i^{(k)}$ are the coefficients corresponding to the eigenvector components.
 177 Since the term $\left\| W_{T_k} - \bar{W}^{(k)} \right\|_2^2$ is a certain constant, the optimization process of weights $W_{T_{k+1}}$ is

178 equivalent to the least upper bound of the $\prod_{t=1+T_k}^{T_{k+1}} (1 - \eta_t \lambda_i^{(k)})$. Thus, the optimization of learning
 179 rate schedule can be formulated as follow,

$$\begin{aligned} & \min_{\eta_{1+T_k}, \eta_{2+T_k}, \dots, \eta_{T_{k+1}}} \max_{\lambda_l^{(k)} \leq \lambda_i^{(k)} \leq \lambda_u^{(k)}} \prod_{t=1+T_k}^{T_{k+1}} \left[(1 - \eta_t \lambda_i^{(k)}) \right]^2 \\ & \text{s.t. } \eta_{1+T_k}, \eta_{2+T_k}, \dots, \eta_{T_{k+1}} \in [\eta_{\min}, \eta_{\max}] \\ & \quad k = 1, 2, \dots, K \end{aligned} \quad (4)$$

180 To solve the constrained min-max problem (4), we adopt an iterative projected gradient method
 181 that alternates between minimizing over the variables η_t and maximizing over the parameters $\lambda_i^{(k)}$.
 182 To avoid repeated numerical optimization, we fit the solutions with a parametric function. Details
 183 regarding the numerical solution and curve fitting process are provided in Appendix B. Then, we
 184 obtain the η_t within the interval $t \in [1 + T_k, T_{k+1}]$ as follows

$$\eta_t = (\eta_{\max} - \eta_{\min}) \frac{2(1 + \cos(\frac{(2(t-T_k)-1)\pi}{2(T_{k+1}-T_k)} + (k-1)\pi))}{2\varphi + (2-\varphi)(1 + \cos(\frac{(2(t-T_k)-1)\pi}{2(T_{k+1}-T_k)} + (k-1)\pi))} + \eta_{\min}, \quad (5)$$

185 where φ is the hyper-parameter controlling the variation speed of the learning rate η_t .

186 When setting ($K > 2$), UBA extends to a multi-phase formulation, which offers three potential
 187 advantages. (i) Hierarchical optimization: by partitioning the budget-aware optimization into mul-
 188 tiple phases, each near a different local minima, it improves approximation accuracy and captures
 189 the dynamic features of the loss surface. (ii) It helps escape saddle points, (iii) It generalizes the
 190 single-phase method, allowing for future extensions. Notably, a single-phase approach ($K = 2$) re-
 191 mains effective, as the robust budget-aware model inherently derives the optimal schedule function,
 192 as shown in Figure 1(a). For consistency with common practice and to better isolate scheduling
 193 effects, this paper focuses primarily on the single-phase implementation. The multi-phase approach
 194 is also compared in Appendix, with its schedule shown in Figure 1(b).

195 We name the proposed schedule Unified Budget-Aware (UBA) schedule for reasons. The robust
 196 budget-aware model minimizes the loss function **uniformly** across the optimization trajectory, re-
 197 sulting in stable performance. UBA provides a reliable, **unified** choice for practitioners, eliminating
 198 the need for case-by-case baseline comparisons and delivering consistent superiority across datasets,
 199 architectures, and training budgets. Lastly, UBA can **uniformly** approximate the behavior of exist-
 200 ing schedules through simple parameter adjustments.

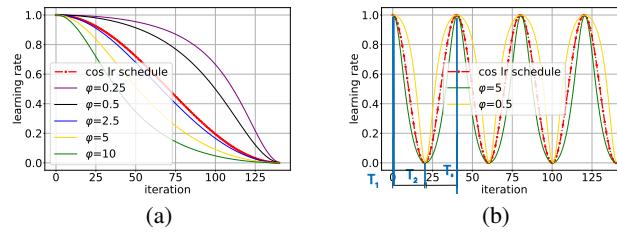


Figure 1: Evolution of the learning rate in UBA schedule across training iterations.

201 3.2 Theoretical analysis

202 **Proposition 1.** *The fit function (5) is the exact closed-form solution to the min-max optimization
 203 problem:*

$$\min_{\eta_{1+T_k}, \eta_{2+T_k}, \dots, \eta_{T_{k+1}}} \max_{\lambda_l^{(k)} \leq \lambda_i^{(k)} \leq \lambda_u^{(k)}} \prod_{t=1+T_k}^{T_{k+1}} \left[\left(1 - \left(\left(\frac{1}{\lambda_l^{(k)}} - \frac{1}{\lambda_u^{(k)}} \right) \eta_t + \frac{1}{\lambda_u^{(k)}} \right) \lambda_i^{(k)} \right) \right]^2 \quad (6)$$

204 when the hyper-parameter φ are determined by $\lambda_l^{(k)}$ and $\lambda_u^{(k)}$ through the relation $\varphi = 2 \frac{\lambda_u^{(k)}}{\lambda_l^{(k)}}$ and
 205 $\eta_{\max} = 1$, $\eta_{\min} = 0$.

206 The min-max model in Proposition 1 represents a special case of our generalized optimization frame-
 207 work. We show that UBA is the exact solution to this special case optimization problem when the
 208 learning rate is scaled by $\left(\left(\frac{1}{\lambda_l^{(k)}} - \frac{1}{\lambda_u^{(k)}} \right) \eta_t + \frac{1}{\lambda_u^{(k)}} \right)$. It provides a theoretical foundation for our
 209 choice of learning rate instead of choosing it heuristically or empirically, adding rigor and justifi-
 210 cation to our approach. Moreover, in this case, φ is linked condition number. It indicates how the
 211 learning rate is shaped by the local curvature of the model, which could guide the optimization pro-
 212 cess more effectively. The relationship between φ and condition number suggests an adaptive nature
 213 for the learning rate. In regions with sharp curvatures (large condition number), φ reduces the learn-
 214 ing rate more rapidly to avoid overshooting. Conversely, in flatter regions, φ allows the learning
 215 rate to remain large for several iterations, facilitating faster convergence. The transformation of the
 216 learning rate trend is shown in Figure 1. *This is a step toward establishing a principled connection*
 217 *between learning rate and local loss landscape geometry along with optimization difficulty.* By this
 218 special case, we generalize that φ is related to optimization difficulty. The empirical results support
 219 these conclusions, with further details in Section 4.3 and Appendix E.5.

220 **Proposition 2.** Consider the training process within the interval $t \in [1+T_k, T_{k+1}]$. When the hyper-
 221 parameter φ is set sufficiently close to 2, the proposed learning rate scheduling formula reduces to
 222 the cosine learning rate schedule.

223 By Proposition 1, the hyper-parameter φ is related to the optimization difficulty. In that case, the
 224 optimization difficulty is explicitly quantified as the condition number and $\varphi = 2\kappa$, where $\kappa = \frac{\lambda_u^{(k)}}{\lambda_l^{(k)}}$.
 225 Furthermore, Proposition 2 shows that when φ approaches 2, the proposed learning rate schedule
 226 converges to the standard cosine schedule. It means that, in regions where the optimization difficulty
 227 is not large (i.e., the curvature is relatively flat), setting $\varphi \approx 2$ allows the learning rate schedule to
 228 naturally reduce to the cosine form. Besides, our schedule can approximate the behavior of existing
 229 schedules (e.g. step decay, cosine annealing, cyclic schedule or Rex schedule) through simple param-
 230 eter adjustments, shown in Table 1. The detail can be found in Appendix C. More importantly, it
 231 outperforms these schedules by training convergence and final accuracy, supporting results can be
 232 found in experiment sections 4.1, 4.2 and Appendix E.4.

Table 1: The adaptive simulation of existing schedules.

Schedule	Parameter adjustments	Schedule	Parameter adjustments
Cosine	$\varphi = 2$	Step	$\varphi = 0, \eta_{\max} = 0.5^k, T_{k+1} - T_k = \text{decaying step}, k = 1, 2, \dots$
Exponential	$\varphi = 30$	Cyclic	$\varphi = 2, k \leftarrow k + 1, T_{k+1} - T_k = \text{cyclic step}, k = 1, 2, \dots$
Rex	$\varphi = 0.8$	OneCycle	$\varphi = 2, k \leftarrow k + 1, T_2 - T_1 = \text{pct_start step}$

233 **Theorem 1.** Let $n_t = H_f^{(k)}(W_t - \bar{W}^{(k)}) - H_f^{(k)}(\xi)(W_t - \bar{W}^{(k)})$ be the stochastic curvature noise
 234 introduced by sampling at iteration t . Assume that the sampling Hessian satisfies: $\mathbb{E}_\xi [n_t n_t^\top] \preceq$
 235 $\sigma^2 H_f^{(k)}$ for some constant σ . Denote $\tau := \frac{4\lambda_l^{(k)}(\eta_{\max} - \eta_{\min})(1 + \cos(\frac{(2(t-T_k)-1)\pi}{2(T_{k+1}-T_k)}))}{(\varphi-2)\pi}$. If we
 236 set the learning rate as the proposed form (5), the loss uncertainty introduced by stochastic gradient
 237 method within the interval $t \in [1+T_k, T_{k+1}]$ can be bounded by two terms,

238 **For** $\varphi > 2$:

$$\begin{aligned}
 & \mathbb{E} [\mathcal{L}(f(W_t, \xi)) - \mathcal{L}(f(\bar{W}^{(k)}, \xi))] \\
 & \leq \left(\frac{4(T_{k+1} - T_k) + (\varphi - 2)\pi}{4(T_{k+1} - T_k) + (\varphi - 2)\pi(t - T_k)} \right)^\tau \cdot \exp(-2\lambda_i^{(k)}\eta_{\min}(t - T_k)) \lambda_u^{(k)} \|W_{1+T_k} - \bar{W}^{(k)}\|_2^2 \\
 & \quad + \sigma^2 \sum_{i=1+T_k}^t \eta_i^2 \sum_{j=1}^N (\lambda_j^{(k)})^2 \exp(-2\lambda_j^{(k)}\eta_{\min}(t - i)) \cdot \left(\frac{4(T_{k+1} - T_k) + (\varphi - 2)\pi(i - T_k)}{4(T_{k+1} - T_k) + (\varphi - 2)\pi(t - T_k)} \right)^\tau
 \end{aligned} \tag{7}$$

239 **For** $\varphi < 2$:

$$\begin{aligned} & \mathbb{E} \left[\mathcal{L}(f(W_t, \xi)) - \mathcal{L}(f(\bar{W}^{(k)}, \xi)) \right] \\ & \leq \left(\frac{(2\varphi + 2\pi - \varphi\pi) - \frac{(2-\varphi)\pi}{(T_{k+1}-T_k)}(t-T_k-0.5)}{(2\varphi + 2\pi - \varphi\pi) + \frac{(2-\varphi)\pi}{2(T_{k+1}-T_k)}} \right)^{-\tau} \cdot \exp \left(-2\lambda_i^{(k)} \eta_{\min}(t-T_k) \right) \lambda_u^{(k)} \|W_{1+T_k} - \bar{W}^{(k)}\|_2^2 \\ & + \sigma^2 \sum_{i=1+T_k}^t \eta_i^2 \sum_{j=1}^N (\lambda_j^{(k)})^2 \exp \left(-2\lambda_j^{(k)} \eta_{\min}(t-i) \right) \cdot \left(\frac{(2\varphi + 2\pi - \varphi\pi) - \frac{(2-\varphi)\pi}{(T_{k+1}-T_k)}(t-T_k-0.5)}{(2\varphi + 2\pi - \varphi\pi) - \frac{(2-\varphi)\pi}{(T_{k+1}-T_k)}(i-T_k-0.5)} \right)^{-\tau} \end{aligned} \quad (8)$$

240 4 Experiment results

241 We conduct comprehensive evaluations along three dimensions: (i) Modality diversity, including
 242 vision and language tasks; (ii) Training budget, including model scales (small to large) and training
 243 iterations (short to long); (iii) Ablation studies, such as parameter sensitivity and cross-optimizer
 244 performance. This evaluation strategy ensures the robustness and generalization of our findings.

245 A key advantage of UBA lies in its critical parameter φ . While an optimal φ can further improve
 246 model performance, we intentionally fix its value across all tasks and architectures to ensure fair
 247 evaluation. We fix $\varphi = 5$ for SGD and $\varphi = 0.5$ for AdamW (see learning rate variations in Figure
 248 1(a)). The rationale for these choices is systematically analyzed in our ablation study 4.3.

249 **Baselines** Research on budgeted-iteration training remains limited, with most existing approaches
 250 focusing on learning rate scheduling strategies [5, 29, 43]. To provide a fair comparison, we adopt
 251 several widely used and empirically effective learning rate schedules as baselines: **Step(SS)** [15],
 252 **Cosine(CS)** [32], **Cyclical (CLR)** and **OneCycle (1C)** [42, 43], **Budgeted training(BT)** [29] and
 253 **Reflected Exponential (REX)** [5]. Details of these baselines are provided in Appendix E.1.

254 4.1 Experiments for Vision Classification Tasks

255 We evaluate our proposed UBA schedule on vision benchmarks (e.g., CIFAR10/100 and ImageNet)
 256 using different architectures (VGG16, ResNet18, ResNet34, ResNet50). In addition, we indepen-
 257 dently train models using fixed epoch budgets of 25%, 50%, and 100% of the maximum training
 258 epochs, without reusing or interpolating results from longer training runs. This setup ensures that
 259 each budget setting is evaluated in isolation, thereby preserving the integrity of comparisons across
 260 low and high training budgets. Table 2 presents the validation accuracy across different training
 261 budgets, comparing UBA against six baseline schedules (see detailed results in Appendix E.3).

262 We find that: (i) UBA demonstrates the strongest performance across all training budgets, achieving
 263 superior results on both small-scale (CIFAR10/100 with ResNet18/34) and large-scale (ImageNet
 264 with ResNet50) benchmarks. This consistent improvement highlights the *generalizability across*
 265 *model and dataset scales*. (ii) UBA outperforms baselines not only at 100% training budget but
 266 also at 25% and 50% iteration budgets, demonstrating its *budget efficiency*, i.e. effectiveness in
 267 computation-constrained scenarios. (iii) Notably, UBA shows robust performance even while the
 268 second-best schedule varies depending on both the datasets, architectures and training budgets. For
 269 practitioners seeking reliable schedules without extensive method selection, UBA provides a default-
 270 strong choice. UBA eliminates the need for case-by-case baseline comparison and delivers *stable*
 271 *superiority and reliability* regardless of datasets, architectures, training budgets and scales.

272 4.2 Experiments for language models

273 We evaluate UBA schedule on the OLMo[18], a truly open language model based on decoder-only
 274 transformer architecture, across diverse benchmarks in language model evaluation. Since large lan-
 275 guage models commonly provide multiple scales to accommodate different compute constraints. We
 276 adjust both model size and training steps to explore budgets. We evaluate UBA on OLMo networks
 277 spanning four parameter scales: 36M, 73M, 151M, and 300M, covering normal to large-scale mod-

Table 2: Validation accuracy for vision classification tasks. We present validation accuracy on vision benchmarks (e.g., CIFAR10/100 and ImageNet) using different architectures (ResNet18, ResNet34, ResNet50) under fixed epoch budgets of 25%, 50%, and 100% of the maximum training epochs.

Schedule	CIFAR10-ResNet18			CIFAR100-ResNet34			ImageNet-ResNet50		
	training budget (epoch(%))			training budget (epoch(%))			training budget (epoch(%))		
	75 (25%)	150 (50%)	300 (100%)	75 (25%)	150 (50%)	300 (100%)	75 (25%)	150 (50%)	300 (100%)
SS	92.41	93.90	93.94	70.85	75.58	77.28	74.84	76.96	78.10
CS	93.66	94.15	95.40	63.88	68.84	70.43	75.99	77.79	79.10
CLR	91.89	92.23	95.32	72.63	73.81	74.80	72.87	74.88	76.91
1C	94.36	95.02	95.48	73.72	75.53	77.90	75.28	77.37	78.79
BT	93.24	94.28	95.55	72.53	75.40	78.49	75.71	77.48	78.66
REX	94.79	94.96	95.59	73.12	75.48	77.99	75.28	77.03	78.46
UBA(ours)	94.54	95.26	95.74	74.57	76.68	78.97	76.00	77.99	79.32

els. Due to space limitation, we defer details such as benchmarks introduction, experimental setting and overall results in Appendix E.4).

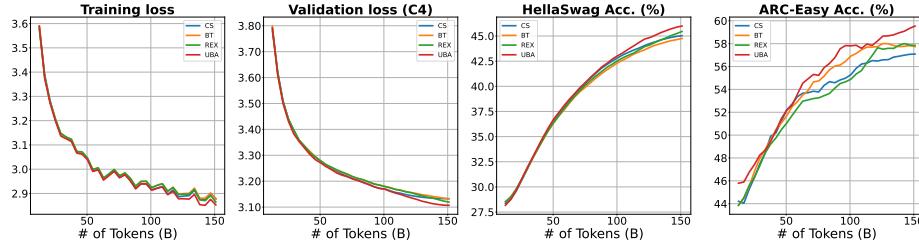


Figure 2: Training dynamics and performance for language tasks under 150B tokens on 300M OLMo. We present the training loss, validation loss, and downstream performance on HSWAG and ARC-E, demonstrating that UBA schedule achieves superior performance.

From Table 3, UBA achieves state-of-the-art performance on approximately 50% of the benchmarks across all scales while the second-best schedule varies. Furthermore, it achieves *consistent superior average performance* among all baselines. It highlights the stable superiority and reliability of UBA, providing a *default-strong choice*. Moreover, it demonstrates significant improvements in SciQ-73M(+1.7) and ARC-E-300M(+2.63), highlighting its ability to enhance generalization across diverse benchmarks. Besides, in Figure 2, UBA consistently achieves lower training loss and validation loss throughout the training, indicating the *efficient training ability* and downstream performance enhancement. Notably, while task-specific tuning of φ can further improve model performance, we intentionally fix its value across all tasks and architectures to ensure fair evaluation. Remarkably, even with this universal φ setting, UBA consistently outperforms baselines on diverse benchmarks, demonstrating *inherent robustness* to varying benchmarks and model scales.

4.3 Ablation Study

Performance across different optimizers UBA originates from the optimization problem under gradient descent dynamics. While modern optimizers (e.g., AdamW [31]) introduce momentum and adaptive mechanisms, they maintain the fundamental property of gradient-based updates. To verify the cross-optimizer performance, we conduct ablation studies between SGD and AdamW optimizers (see detailed results in Appendix E.5). The results show that UBA achieves SOTAs on both SGD and AdamW optimizers, demonstrating the *cross-optimizer robustness* of UBA. This highlights its *broad applicability* beyond standard gradient descent.

Parameter analysis of φ We perform a sensitivity analysis of $\varphi \in \{0.25, 0.5, 1.0, 2.5, 5, 10\}$ in equation (5), which controls the variation speed of learning rate (see detailed setting, results and analyses in Appendix E.6). Experiments 11 show AdamW prefers smaller φ while SGD requires larger φ . We attribute this phenomenon to the preconditioning effect of AdamW. As the relationship

Table 3: Performance comparison between UBA and the best-performing baseline schedules on OLMo. We report the results of top-performing baseline schedule for the corresponding benchmark and model scale. The names of the top-performing baseline schedules are listed below the results. Bold values indicate UBA’s superiority (See overall results of each schedule in Appendix E.4).

Size	Sched.	Benchmark(accuracy %)										
		PIQA	HSWAG	OBQA	SciQ	ARC-E	ARC-C	COPA	SIQA	SOC	OTH	Avg.
36M	Best baseline (CLR) (1C)	61.15	27.91	27.80	68.10	45.26	23.41	63.00	41.45	25.18	29.86	40.92
	UBA	60.39	27.98	27.40	68.20	45.79	21.74	63.00	40.69	24.33	30.14	40.97
73M	Best baseline (REX) (REX)	62.46	30.22	28.80	72.60	47.54	26.42	66.00	41.97	26.32	29.34	42.53
	UBA	63.17	30.09	28.80	74.30	45.79	22.74	65.00	41.45	27.13	28.85	42.73
150M	Best baseline (CS) (REX)	66.00	35.72	32.80	78.20	53.16	26.42	69.00	43.71	27.68	32.64	45.91
	UBA	65.23	35.50	29.80	78.30	50.35	27.42	67.00	43.50	29.29	32.72	45.91
300M	Best baseline (REX) (REX)	70.18	46.30	33.40	84.40	57.72	28.43	72.00	44.58	29.50	36.74	49.70
	UBA	69.48	46.44	34.60	83.90	60.35	29.10	72.00	44.42	28.53	35.41	50.42

303 formalized in Proposition 1), smaller φ is favored for low optimization difficult, while larger φ
 304 is beneficial for difficult optimization. AdamW adapts the learning rate by scaling gradients with
 305 the second moment estimate $\sqrt{v_t}$, implicitly reducing the condition number of the optimization
 306 landscape. Thus the schedule with smaller φ is preferred. In contrast, SGD lacks such adaptive
 307 mechanisms, thus a larger φ is more suitable for this situation. This alignment between theory and
 308 experiment underscores the importance of tailoring φ to the optimizer’s characteristics. *Overall, φ*
 309 *is related to a generalized optimization difficulty, where optimization precondition effect, datasets*
 310 *distribution and network architecture are all related to optimization difficulty.*

311 **Performance across different periods** Our schedule has a periodic phase-based learning rate ad-
 312 justment setting, where the learning rate at the k -th phase is dynamically determined by the k -th
 313 local minimum of the loss landscape. To validate our scheduling strategy, we conduct experiments
 314 by varying K (see detailed results in Appendix E.7). The results suggests that multi-phase schedul-
 315 ing captures the dynamic features of the loss surface more finely, but it needs careful selection for
 316 φ , where φ reflects optimization difficulty. However, selecting optimal φ values per phase for multi-
 317 phase remains non-trivial, which motivates future work on automated landscape-aware φ tuning.

318 5 Conclusion

319 In this paper, we construct a unified budget-aware training optimization framework that inherently
 320 adapts to landscape curvature variations, enhancing training robustness. From this optimization
 321 framework, we propose the Unified Budget-Aware (UBA) schedule. Extensive experiments demon-
 322 strate UBA *consistently surpasses* the commonly-used schedules across diverse vision and language
 323 tasks, spanning network architectures (e.g., ResNet, OLMo) and scales, under different training-
 324 iteration budgets.

325 Theoretically and empirically, we observe that the parameter φ correlates with the optimization dif-
 326 ficulty of the training process, influences the optimal choice of φ and UBA’s performance. However,
 327 the explicit relationship between φ and optimization difficulty remains unexplored and no estab-
 328 lished evaluation metric exists to quantify the optimization difficulty. These limitations motivate
 329 future work on optimization difficulty-aware φ tuning.

330 Despite these open questions, the effectiveness, implementation simplicity and ease of tuning make
 331 the UBA a practical, must-try schedule for deep learning practitioners. We hope this work could
 332 motivate more studies on learning rate scheduling.

333 **References**

- 334 [1] Ebtesam Almazrouei, Hamza Alobeidli, Abdulaziz Alshamsi, Alessandro Cappelli, Ruxandra Cojocaru,
335 M  rouane Debbah,  tienne Goffinet, Daniel Hesslow, Julien Launay, Quentin Malartic, et al. The falcon
336 series of open language models. *arXiv preprint arXiv:2311.16867*, 2023.
- 337 [2] Eric Arazo, Diego Ortego, Paul Albert, Noel E O'Connor, and Kevin McGuinness. How important is
338 importance sampling for deep budgeted training? *arXiv preprint arXiv:2110.14283*, 2021.
- 339 [3] Yonatan Bisk, Rowan Zellers, Jianfeng Gao, Yejin Choi, et al. Piqa: Reasoning about physical common-
340 sense in natural language. In *Proceedings of the AAAI conference on artificial intelligence*, volume 34,
341 pages 7432–7439, 2020.
- 342 [4] Tom B Brown. Language models are few-shot learners. In *Advances in Neural Information Processing
343 Systems*, pages 1877–1901, 2020.
- 344 [5] John Chen, Cameron Wolfe, and Tasos Kyrillidis. Rex: Revisiting budgeted training with an improved
345 schedule. *Proceedings of Machine Learning and Systems*, 4:64–76, 2022.
- 346 [6] Tianlong Chen, Yu Cheng, Zhe Gan, Jingjing Liu, and Zhangyang Wang. Data-efficient gan training
347 beyond (just) augmentations: A lottery ticket perspective. *Advances in Neural Information Processing
348 Systems*, 34:20941–20955, 2021.
- 349 [7] Xiangning Chen, Chen Liang, Da Huang, Esteban Real, Kaiyuan Wang, Hieu Pham, Xuanyi Dong, Thang
350 Luong, Cho-Jui Hsieh, Yifeng Lu, et al. Symbolic discovery of optimization algorithms. *Advances in
351 neural information processing systems*, 36, 2024.
- 352 [8] Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam Roberts,
353 Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, et al. Palm: Scaling language
354 modeling with pathways. *Journal of Machine Learning Research*, 24(240):1–113, 2023.
- 355 [9] Christopher Clark, Kenton Lee, Ming-Wei Chang, Tom Kwiatkowski, Michael Collins, and Kristina
356 Toutanova. Boolq: Exploring the surprising difficulty of natural yes/no questions. *arXiv preprint
357 arXiv:1905.10044*, 2019.
- 358 [10] Peter Clark, Isaac Cowhey, Oren Etzioni, Tushar Khot, Ashish Sabharwal, Carissa Schoenick, and Oyvind
359 Tafjord. Think you have solved question answering? try arc, the ai2 reasoning challenge. *arXiv preprint
360 arXiv:1803.05457*, 2018.
- 361 [11] Simon Du, Jason Lee, Yuandong Tian, Aarti Singh, and Barnabas Poczos. Gradient descent learns one-
362 hidden-layer cnn: Dont be afraid of spurious local minima. In *International Conference on Machine
363 Learning*, pages 1339–1348. PMLR, 2018.
- 364 [12] John Duchi, Elad Hazan, and Yoram Singer. Adaptive subgradient methods for online learning and
365 stochastic optimization. *Journal of machine learning research*, 12(7), 2011.
- 366 [13] Donald A. Flanders and George Shortley. Numerical determination of fundamental modes. *Journal of
367 Applied Physics*, 21(12):1326–1332, 1951.
- 368 [14] Leo Gao, Jonathan Tow, Baber Abbasi, Stella Biderman, Sid Black, Anthony DiPofi, Charles Foster,
369 Laurence Golding, and Jeffrey Hsu. A framework for few-shot language model evaluation. December
370 2023. URL <https://zenodo.org/records/10256836>.
- 371 [15] Rong Ge, Sham M Kakade, Rahul Kidambi, and Praneeth Netrapalli. The step decay schedule: A near
372 optimal, geometrically decaying learning rate procedure for least squares. In *Advances in Neural Infor-
373 mation Processing Systems*, volume 32, 2019.
- 374 [16] Jonas Geiping and Tom Goldstein. Cramming: Training a language model on a single gpu in one day. In
375 *International Conference on Machine Learning*, pages 11117–11143. PMLR, 2023.
- 376 [17] Andrew Gordon, Zornitsa Kozareva, and Melissa Roemmele. Semeval-2012 task 7: Choice of plausible
377 alternatives: An evaluation of commonsense causal reasoning. In **SEM 2012: The First Joint Conference
378 on Lexical and Computational Semantics–Volume 1: Proceedings of the main conference and the shared
379 task, and Volume 2: Proceedings of the Sixth International Workshop on Semantic Evaluation (SemEval
380 2012)*, pages 394–398, 2012.
- 381 [18] Dirk Groeneveld, Iz Beltagy, Pete Walsh, Akshita Bhagia, Rodney Kinney, Oyvind Tafjord, Ananya Harsh
382 Jha, Hamish Ivison, Ian Magnusson, Yizhong Wang, et al. Olmo: Accelerating the science of language
383 models. *arXiv preprint arXiv:2402.00838*, 2024.

- 384 [19] Alex Hägele, Elie Bakouch, Atli Kosson, Leandro Von Werra, Martin Jaggi, et al. Scaling laws and
 385 compute-optimal training beyond fixed training durations. *Advances in Neural Information Processing*
 386 *Systems*, 37:76232–76264, 2024.
- 387 [20] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition.
 388 In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016.
- 389 [21] Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Stein-
 390 hardt. Measuring massive multitask language understanding. *arXiv preprint arXiv:2009.03300*, 2020.
- 391 [22] Geoffrey Hinton, Nitish Srivastava, and Kevin Swersky. Neural networks for machine learning lecture 6a
 392 overview of mini-batch gradient descent. *Cited on*, 14(8):2, 2012.
- 393 [23] Jordan Hoffmann, Sebastian Borgeaud, Arthur Mensch, Elena Buchatskaya, Trevor Cai, Eliza Rutherford,
 394 Diego de Las Casas, Lisa Anne Hendricks, Johannes Welbl, Aidan Clark, et al. Training compute-optimal
 395 large language models. *arXiv preprint arXiv:2203.15556*, 2022.
- 396 [24] Shengding Hu, Yuge Tu, Xu Han, Chaoqun He, Ganqu Cui, Xiang Long, Zhi Zheng, Yewei Fang, Yuxiang
 397 Huang, Weilin Zhao, et al. Minicpm: Unveiling the potential of small language models with scalable
 398 training strategies. *arXiv preprint arXiv:2404.06395*, 2024.
- 399 [25] Peter Izsak, Moshe Berchansky, and Omer Levy. How to train bert with an academic budget. *arXiv*
 400 *preprint arXiv:2104.07705*, 2021.
- 401 [26] Jean Kaddour. Stop wasting my time! saving days of imagenet and bert training with latest weight
 402 averaging. *arXiv preprint arXiv:2209.14981*, 2022.
- 403 [27] Krishnateja Killamsetty, Sivasubramanian Durga, Ganesh Ramakrishnan, Abir De, and Rishabh Iyer.
 404 Grad-match: Gradient matching based data subset selection for efficient deep model training. In *Interna-*
 405 *tional Conference on Machine Learning*, pages 5464–5474. PMLR, 2021.
- 406 [28] Diederik P Kingma. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- 407 [29] Mengtian Li, Ersin Yumer, and Deva Ramanan. Budgeted training: Rethinking deep neural network
 408 training under resource constraints. *arXiv preprint arXiv:1905.04753*, 2019.
- 409 [30] Zhuohan Li, Eric Wallace, Sheng Shen, Kevin Lin, Kurt Keutzer, Dan Klein, and Joey Gonzalez. Train big,
 410 then compress: Rethinking model size for efficient training and inference of transformers. In *International*
 411 *Conference on machine learning*, pages 5958–5968. PMLR, 2020.
- 412 [31] I Loshchilov. Decoupled weight decay regularization. *arXiv preprint arXiv:1711.05101*, 2017.
- 413 [32] Ilya Loshchilov and Frank Hutter. Sgdr: Stochastic gradient descent with warm restarts. *arXiv preprint*
 414 *arXiv:1608.03983*, 2016.
- 415 [33] Kairong Luo, Haodong Wen, Shengding Hu, Zhenbo Sun, Zhiyuan Liu, Maosong Sun, Kaifeng Lyu,
 416 and Wenguang Chen. A multi-power law for loss curve prediction across learning rate schedules. In *International*
 417 *Conference on Learning Representations (ICLR)*, 2025.
- 418 [34] Todor Mihaylov, Peter Clark, Tushar Khot, and Ashish Sabharwal. Can a suit of armor conduct electricity?
 419 a new dataset for open book question answering. *arXiv preprint arXiv:1809.02789*, 2018.
- 420 [35] Rui Pan, Haishan Ye, and Tong Zhang. Eigencurve: Optimal learning rate schedule for sgd on quadratic
 421 objectives with skewed hessian spectrums. *arXiv preprint arXiv:2110.14109*, 2021.
- 422 [36] Rui Pan, Shizhe Diao, Jianlin Chen, and Tong Zhang. Extremebert: A toolkit for accelerating pretraining
 423 of customized bert. *arXiv preprint arXiv:2211.17201*, 2022.
- 424 [37] Xuran Pan, Xuan Jin, Yuan He, Shiji Song, Gao Huang, et al. Budgeted training for vision transformer.
 425 In *The Eleventh International Conference on Learning Representations*, 2022.
- 426 [38] Keisuke Sakaguchi, Ronan Le Bras, Chandra Bhagavatula, and Yejin Choi. Winogrande: An adversarial
 427 winograd schema challenge at scale. *Communications of the ACM*, 64(9):99–106, 2021.
- 428 [39] Maarten Sap, Hannah Rashkin, Derek Chen, Ronan LeBras, and Yejin Choi. Socialiq: Commonsense
 429 reasoning about social interactions. *arXiv preprint arXiv:1904.09728*, 2019.
- 430 [40] Li Shen, Yan Sun, Zhiyuan Yu, Liang Ding, Xinmei Tian, and Dacheng Tao. On efficient training of
 431 large-scale deep learning models: A literature review. *arXiv preprint arXiv:2304.03589*, 2023.

- 432 [41] Yong Shi, Anda Tang, Lingfeng Niu, and Ruizhi Zhou. Sparse optimization guided pruning for neural
433 networks. *Neurocomputing*, 574:127280, 2024.
- 434 [42] Leslie N Smith. Cyclical learning rates for training neural networks. In *2017 IEEE winter conference on*
435 *applications of computer vision (WACV)*, pages 464–472. IEEE, 2017.
- 436 [43] Leslie N Smith and Nicholay Topin. Super-convergence: Very fast training of neural networks using large
437 learning rates. In *Artificial intelligence and machine learning for multi-domain operations applications*,
438 volume 11006, pages 369–386. SPIE, 2019.
- 439 [44] Vivienne Sze, Yu-Hsin Chen, Joel Emer, Amr Suleiman, and Zhengdong Zhang. Hardware for machine
440 learning: Challenges and opportunities. In *2017 IEEE custom integrated circuits conference (CICC)*,
441 pages 1–8. IEEE, 2017.
- 442 [45] Alon Talmor, Jonathan Herzig, Nicholas Lourie, and Jonathan Berant. Commonsenseqa: A question
443 answering challenge targeting commonsense knowledge. *arXiv preprint arXiv:1811.00937*, 2018.
- 444 [46] Hugo Touvron, Matthieu Cord, Matthijs Douze, Francisco Massa, Alexandre Sablayrolles, and Hervé
445 Jégou. Training data-efficient image transformers & distillation through attention. In *International con-*
446 *ference on machine learning*, pages 10347–10357. PMLR, 2021.
- 447 [47] Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay
448 Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open foundation and
449 fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023.
- 450 [48] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz
451 Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing*
452 *systems*, 30, 2017.
- 453 [49] Johannes Welbl, Nelson F Liu, and Matt Gardner. Crowdsourcing multiple choice science questions.
454 *arXiv preprint arXiv:1707.06209*, 2017.
- 455 [50] Ashia C Wilson, Rebecca Roelofs, Mitchell Stern, Nati Srebro, and Benjamin Recht. The marginal value
456 of adaptive gradient methods in machine learning. *Advances in neural information processing*
457 *systems*, 30, 2017.
- 458 [51] Xingyu Xie, Pan Zhou, Huan Li, Zhouchen Lin, and Shuicheng Yan. Adan: Adaptive nesterov momen-
459 tumb algorithm for faster optimizing deep models. *IEEE Transactions on Pattern Analysis and Machine*
460 *Intelligence*, 2024.
- 461 [52] Kh Shahriya Zaman, Mamun Bin Ibne Reaz, Sawal Hamid Md Ali, Ahmad Ashrif A Bakar, and Muham-
462 mad Enamul Hoque Chowdhury. Custom hardware architectures for deep learning on portable devices: a
463 review. *IEEE Transactions on Neural Networks and Learning Systems*, 33(11):6068–6088, 2021.
- 464 [53] Matthew D Zeiler. Adadelta: an adaptive learning rate method. *arXiv preprint arXiv:1212.5701*, 2012.
- 465 [54] Rowan Zellers, Ari Holtzman, Yonatan Bisk, Ali Farhadi, and Yejin Choi. Hellaswag: Can a machine
466 really finish your sentence? *arXiv preprint arXiv:1905.07830*, 2019.
- 467 [55] Jian Zhang, Lei Qi, Yinghuan Shi, and Yang Gao. Exploring flat minima for domain generalization with
468 large learning rates. *IEEE Transactions on Knowledge and Data Engineering*, 2024.
- 469 [56] Zhiyuan Zhang, Ruixuan Luo, Qi Su, and Xu Sun. Ga-sam: Gradient-strength based adaptive sharpness-
470 aware minimization for improved generalization. *arXiv preprint arXiv:2210.06895*, 2022.
- 471 [57] Xunyu Zhu, Jian Li, Yong Liu, Can Ma, and Weiping Wang. A survey on model compression for large
472 language models. *Transactions of the Association for Computational Linguistics*, 12:1556–1577, 2024.

473 **NeurIPS Paper Checklist**

474 **1. Claims**

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

477 Answer: [Yes]

478 Justification: The main claims in the abstract and introduction (e.g., UBA consistently sur-
479 passes the commonly-used schedules on vision and language tasks and sets new SoTAs
480 for many networks and frameworks, e.g. ResNet and OLMo, under different training iter-
481 ations.) are fully aligned with the papers contributions and scope. We provide an in-depth
482 experiments and analyses of the proposed learning schedule across vision and language
483 tasks, networks, scales and training budgets.

484 Guidelines:

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

494 **2. Limitations**

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

496 Answer: [Yes]

497 Justification: We explicitly discuss limitations and future work in Section 5. Besides, the
498 assumptions can be found in Appendix A. We test our method across vision and language
499 tasks, various scales of datasets and networks, under different iteration budgets.

500 Guidelines:

- 501 • The answer NA means that the paper has no limitation while the answer No means
502 that the paper has limitations, but those are not discussed in the paper.
- 503 • The authors are encouraged to create a separate "Limitations" section in their paper.
- 504 • The paper should point out any strong assumptions and how robust the results are to
505 violations of these assumptions (e.g., independence assumptions, noiseless settings,
506 model well-specification, asymptotic approximations only holding locally). The au-
507 thors should reflect on how these assumptions might be violated in practice and what
508 the implications would be.
- 509 • The authors should reflect on the scope of the claims made, e.g., if the approach was
510 only tested on a few datasets or with a few runs. In general, empirical results often
511 depend on implicit assumptions, which should be articulated.
- 512 • The authors should reflect on the factors that influence the performance of the ap-
513 proach. For example, a facial recognition algorithm may perform poorly when image
514 resolution is low or images are taken in low lighting. Or a speech-to-text system might
515 not be used reliably to provide closed captions for online lectures because it fails to
516 handle technical jargon.
- 517 • The authors should discuss the computational efficiency of the proposed algorithms
518 and how they scale with dataset size.
- 519 • If applicable, the authors should discuss possible limitations of their approach to ad-
520 dress problems of privacy and fairness.
- 521 • While the authors might fear that complete honesty about limitations might be used by
522 reviewers as grounds for rejection, a worse outcome might be that reviewers discover
523 limitations that aren't acknowledged in the paper. The authors should use their best
524 judgment and recognize that individual actions in favor of transparency play an impor-
525 tant role in developing norms that preserve the integrity of the community. Reviewers
526 will be specifically instructed to not penalize honesty concerning limitations.

527 **3. Theory assumptions and proofs**

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

530 Answer: [Yes]

531 Justification: The full set of assumptions are in Appendix A. The complete derivations and
532 proof are in Appendix A and D.

533 Guidelines:

- 534 • The answer NA means that the paper does not include theoretical results.
- 535 • All the theorems, formulas, and proofs in the paper should be numbered and cross-
536 referenced.
- 537 • All assumptions should be clearly stated or referenced in the statement of any theo-
538 rems.
- 539 • The proofs can either appear in the main paper or the supplemental material, but if
540 they appear in the supplemental material, the authors are encouraged to provide a
541 short proof sketch to provide intuition.
- 542 • Inversely, any informal proof provided in the core of the paper should be comple-
543 mented by formal proofs provided in appendix or supplemental material.
- 544 • Theorems and Lemmas that the proof relies upon should be properly referenced.

545 **4. Experimental result reproducibility**

546 Question: Does the paper fully disclose all the information needed to reproduce the main
547 experimental results of the paper to the extent that it affects the main claims and/or conclu-
548 sions of the paper (regardless of whether the code and data are provided or not)?

549 Answer: [Yes]

550 Justification: In Appendix E, we carefully describe provide complete information needed
551 to reproduce the experimental results, including hyper-parameters(initial learning rates,
552 batch sizes, configurations of optimizer and schedules) for all tasks, implementation ver-
553 sions(such as PyTorch 2.4.1+cu118 on CIFAR), hardware configurations(GPU models,
554 memory), full results of ablation studies (see Appendix E). Code and data are also available.

555 Guidelines:

- 556 • The answer NA means that the paper does not include experiments.
- 557 • If the paper includes experiments, a No answer to this question will not be perceived
558 well by the reviewers: Making the paper reproducible is important, regardless of
559 whether the code and data are provided or not.
- 560 • If the contribution is a dataset and/or model, the authors should describe the steps
561 taken to make their results reproducible or verifiable.
- 562 • Depending on the contribution, reproducibility can be accomplished in various ways.
563 For example, if the contribution is a novel architecture, describing the architecture
564 fully might suffice, or if the contribution is a specific model and empirical evaluation,
565 it may be necessary to either make it possible for others to replicate the model with
566 the same dataset, or provide access to the model. In general, releasing code and data
567 is often one good way to accomplish this, but reproducibility can also be provided via
568 detailed instructions for how to replicate the results, access to a hosted model (e.g., in
569 the case of a large language model), releasing of a model checkpoint, or other means
570 that are appropriate to the research performed.
- 571 • While NeurIPS does not require releasing code, the conference does require all sub-
572 missions to provide some reasonable avenue for reproducibility, which may depend
573 on the nature of the contribution. For example
 - 574 (a) If the contribution is primarily a new algorithm, the paper should make it clear
575 how to reproduce that algorithm.
 - 576 (b) If the contribution is primarily a new model architecture, the paper should describe
577 the architecture clearly and fully.
 - 578 (c) If the contribution is a new model (e.g., a large language model), then there should
579 either be a way to access this model for reproducing the results or a way to re-
580 produce the model (e.g., with an open-source dataset or instructions for how to
581 construct the dataset).

582 (d) We recognize that reproducibility may be tricky in some cases, in which case au-
583 thors are welcome to describe the particular way they provide for reproducibility.
584 In the case of closed-source models, it may be that access to the model is limited in
585 some way (e.g., to registered users), but it should be possible for other researchers
586 to have some path to reproducing or verifying the results.

587 **5. Open access to data and code**

588 Question: Does the paper provide open access to the data and code, with sufficient instruc-
589 tions to faithfully reproduce the main experimental results, as described in supplemental
590 material?

591 Answer: [Yes]

592 Justification: We provide open access to the code in the Appendix. Moreover, we provide
593 sufficient instructions to faithfully reproduce the main experimental results in Appendix E.

594 Guidelines:

- 595 • The answer NA means that paper does not include experiments requiring code.
- 596 • Please see the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- 597 • While we encourage the release of code and data, we understand that this might not
598 be possible, so No is an acceptable answer. Papers cannot be rejected simply for not
600 including code, unless this is central to the contribution (e.g., for a new open-source
601 benchmark).
- 602 • The instructions should contain the exact command and environment needed to run to
603 reproduce the results. See the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- 604 • The authors should provide instructions on data access and preparation, including how
605 to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- 606 • The authors should provide scripts to reproduce all experimental results for the new
607 proposed method and baselines. If only a subset of experiments are reproducible, they
608 should state which ones are omitted from the script and why.
- 609 • At submission time, to preserve anonymity, the authors should release anonymized
610 versions (if applicable).
- 611 • Providing as much information as possible in supplemental material (appended to the
612 paper) is recommended, but including URLs to data and code is permitted.

614 **6. Experimental setting/details**

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

618 Answer: [Yes]

619 Justification: These details are described in Appendix E.

620 Guidelines:

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

626 **7. Experiment statistical significance**

627 Question: Does the paper report error bars suitably and correctly defined or other appropri-
628 ate information about the statistical significance of the experiments?

629 Answer: [No]

630 Justification: Due to the high computational cost of training large-scale models (e.g., Trans-
631 formers on ImageNet), we report results from a single run with fixed random seeds. To
632 ensure reliability, we list experimental setting and provide code for deterministic reproduc-
633 tion. To maintain uniform evaluation protocols across all experiments (large/small-scale),

634 we prioritized identical statistical reporting standards. This avoids selective rigor that could
635 mislead comparisons. Besides, we conduct parameter analyses to improve reliability. More-
636 over, small-scale results align with the expectations and results on large-scale experiments,
637 reducing the need for empirical variance quantification. However, we provide error bars of
638 curve fitting on numerical solution.

639 Guidelines:

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

660 8. Experiments compute resources

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

664 Answer: [Yes]

665 Justification: The computation resource is described in Appendix E. we carefully describe
666 provide complete information including type of compute workers and memory.

667 Guidelines:

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

676 9. Code of ethics

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

679 Answer: [Yes]

680 Justification: We conform with the NeurIPS code of Ethics.

681 Guidelines:

- 682 • The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- 683 • If the authors answer No, they should explain the special circumstances that require a
684 deviation from the Code of Ethics.

- 685 • The authors should make sure to preserve anonymity (e.g., if there is a special consid-
686 eration due to laws or regulations in their jurisdiction).

687 **10. Broader impacts**

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

690 Answer: [Yes]

691 Justification: The paper discusses societal impacts in Section 5. It highlights potential ben-
692 efits, such as improving model performance and offering a practical scheduling strategy for
693 practitioners, which may inspire further research on learning rate scheduling. Since we note
694 that the proposed method is a generic optimization algorithm with no explicit negative soci-
695 etal impacts, indirect risks such as computational costs from sub-optimal hyper-parameter
696 selection are not thoroughly examined. We provide the hyper-parameter selection guideline
697 through analyses to mitigate this indirect risk.

698 Guidelines:

- 699 • The answer NA means that there is no societal impact of the work performed.
700 • If the authors answer NA or No, they should explain why their work has no societal
701 impact or why the paper does not address societal impact.
702 • Examples of negative societal impacts include potential malicious or unintended uses
703 (e.g., disinformation, generating fake profiles, surveillance), fairness considerations
704 (e.g., deployment of technologies that could make decisions that unfairly impact spe-
705 cific groups), privacy considerations, and security considerations.
706 • The conference expects that many papers will be foundational research and not tied
707 to particular applications, let alone deployments. However, if there is a direct path to
708 any negative applications, the authors should point it out. For example, it is legitimate
709 to point out that an improvement in the quality of generative models could be used to
710 generate deepfakes for disinformation. On the other hand, it is not needed to point out
711 that a generic algorithm for optimizing neural networks could enable people to train
712 models that generate Deepfakes faster.
713 • The authors should consider possible harms that could arise when the technology is
714 being used as intended and functioning correctly, harms that could arise when the
715 technology is being used as intended but gives incorrect results, and harms following
716 from (intentional or unintentional) misuse of the technology.
717 • If there are negative societal impacts, the authors could also discuss possible mitiga-
718 tion strategies (e.g., gated release of models, providing defenses in addition to attacks,
719 mechanisms for monitoring misuse, mechanisms to monitor how a system learns from
720 feedback over time, improving the efficiency and accessibility of ML).

721 **11. Safeguards**

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

725 Answer: [NA]

726 Justification: Our paper poses no such risks.

727 Guidelines:

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

738 **12. Licenses for existing assets**

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

742 Answer: [Yes]

743 Justification: All assets are properly cited.

744 Guidelines:

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

760 **13. New assets**

761 Question: Are new assets introduced in the paper well documented and is the documenta-
762 tion provided alongside the assets?

763 Answer: [NA]

764 Justification: The paper does not release new assets.

765 Guidelines:

- 766 • The answer NA means that the paper does not release new assets.
- 767 • Researchers should communicate the details of the dataset/code/model as part of their
768 submissions via structured templates. This includes details about training, license,
769 limitations, etc.
- 770 • The paper should discuss whether and how consent was obtained from people whose
771 asset is used.
- 772 • At submission time, remember to anonymize your assets (if applicable). You can
773 either create an anonymized URL or include an anonymized zip file.

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

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

778 Answer: [NA]

779 Justification: The paper does not involve crowdsourcing nor research with human subjects.

780 Guidelines:

- 781 • The answer NA means that the paper does not involve crowdsourcing nor research
782 with human subjects.
- 783 • Including this information in the supplemental material is fine, but if the main contri-
784 bution of the paper involves human subjects, then as much detail as possible should
785 be included in the main paper.
- 786 • According to the NeurIPS Code of Ethics, workers involved in data collection, cura-
787 tion, or other labor should be paid at least the minimum wage in the country of the
788 data collector.

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

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

795 Answer: [NA]

796 Justification: The paper does not involve crowdsourcing nor research with human subjects.

797 Guidelines:

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

808 **16. Declaration of LLM usage**

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

813 Answer: [NA]

814 Justification: The core method development in this research does not involve LLMs as any
815 important, original, or non-standard components.

816 Guidelines:

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