
On Your Mark, Get Set, Warmup!

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 It is common in deep learning to warm up the learning rate η , often by a linear
2 schedule between $\eta_{\text{init}} = 0$ and a predetermined target η_{trgt} . In this paper, we
3 show through systematic experiments using SGD and Adam that the overwhelming
4 benefit of warmup arises from allowing the network to tolerate larger η_{trgt} by forcing
5 the network to more well-conditioned areas of the loss landscape. The ability to
6 handle larger η_{trgt} makes hyperparameter tuning more robust while improving
7 the final performance. We uncover different regimes of operation during the
8 warmup period, depending on whether training starts off in a progressive sharpening
9 or sharpness reduction phase, which in turn depends on the initialization and
10 parameterization. We also suggest an initialization for the variance in Adam which
11 provides benefits similar to warmup.

12 1 Introduction

13 One of the most important choices to make in gradient-based optimization is the learning rate (step
14 size) η . If η is too small, then learning may take place too slowly or the model might get stuck in
15 unfavorable regions of the loss landscape. If η is too large, training will typically diverge. In practice,
16 it is common to pick a dynamical learning rate schedule η_t [2, 4, 39, 26]. Modern learning rate
17 schedules for deep learning typically consist of a warmup period where η_t is increased linearly from
18 zero to a target value η_{trgt} over a warmup time T_{wrm} [13, 33]. After the warmup period, it is common
19 to eventually decay the learning rate, for example via a cosine decay schedule [33, 26, 39].

20 Given that warmup is standard in the practitioner’s toolkit, it is important to understand it deeply and
21 identify improvements. In modern settings, perhaps the earliest work to use warmup was [14], which
22 used a small constant learning rate for the first few epochs of training and then switched to a larger
23 learning rate. A linear warmup schedule was later introduced in [13]. The intuition given was that to
24 scale the minibatch size in SGD by a factor of k , it is natural to also scale the learning rate by a factor
25 of k , provided the model is not changing too rapidly and successive gradients are roughly aligned.
26 However at the beginning of training, the model is changing rapidly, so it is natural to start with a
27 lower learning rate and gradually increase it to the target value after the network has stabilized. Other
28 explanations suggest that since the network is initialized randomly, the gradient steps at the beginning
29 of training are not meaningful, and thus it would be harmful to take large steps in such directions [39],
30 so it makes sense to take smaller steps early in training. The analysis by [12] suggests that warmup
31 primarily limits the magnitude of weight updates in the deeper layers, preventing large instabilities. It
32 has also been suggested that the key benefit of warmup arises for adaptive optimizers, such as Adam:
33 [23] argues that the variance of the adaptive learning rate is large during early training because the
34 network has seen too few training samples; it is asserted that this large variance is harmful, and that
35 warmup acts as a variance reduction method by allowing the network to collect accurate statistics of
36 the gradient moments before using larger learning rates. Alternatively, it is also sometimes stated that
37 the initialization may start the model off at places in parameter space that are unstable, difficult to
38 optimize, and easily lead to divergence, and that warmup can help alleviate this [39].

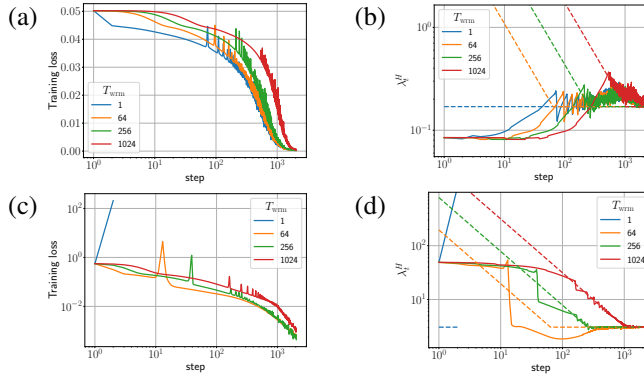


Figure 1: Training loss and sharpness trajectories of FCNs trained on a 5k subset of CIFAR-10 with MSE loss using GD. The dashed lines in the sharpness figures illustrate the instability thresholds $2/\eta_t$. (top) μP with $\eta_{\text{trgt}} = 1/\lambda_0^H$, (bottom) SP with $\eta_{\text{trgt}} = 32/\lambda_0^H$. Similar mechanisms are observed across different architectures, loss functions, and mini-batch sizes, as shown in Appendix E.

39 The above explanations are varied and do not clearly demonstrate why and to what extent warmup is
 40 necessary. A loss landscape perspective was given in [10] (and summarized in [26] Ch. 8), which
 41 argued that an important effect of warmup is to gradually reduce the sharpness (the top eigenvalue
 42 of the Hessian of the loss), thus causing the model to leave poorly conditioned areas of the loss
 43 landscape and move towards flatter regions which can tolerate larger learning rates. They argue that
 44 the mechanism for this is similar to the dynamical stability (catapult) mechanisms studied in [34, 22].

45 **Our contributions.** In this paper, we perform extensive studies on the effect of learning rate
 46 warmup across a variety of architectures, initializations and parameterizations, datasets, and for
 47 both SGD and Adam optimizers. We demonstrate through systematic experiments that by far the
 48 primary benefit of learning rate warmup is to allow the network to tolerate larger learning rates than
 49 it otherwise would have. This builds on the observations of [10] by showing that any other benefits
 50 are marginal, disentangling the effect of warmup duration and target learning rate, and by extending
 51 the empirical evidence to include adaptive optimizers and Transformers.

52 2 Notations and Preliminaries

53 **Sharpness:** The sharpness is defined as the maximum eigenvalue of the Hessian of the loss
 54 $\lambda_t^H := \lambda_{\max}(\nabla_{\theta}^2 L)$ at training step t . For adaptive optimizers with pre-conditioner P , $\lambda^{P^{-1}H} :=$
 55 $\lambda_{\max}(P^{-1}\nabla_{\theta}^2 L)$ denotes the pre-conditioned sharpness. For details on Adam’s pre-conditioner, see
 56 Appendix D.3.

57 **Linear Warmup:** This is defined by the schedule $\eta_t = \eta_{\text{init}} + (\eta_{\text{trgt}} - \eta_{\text{init}}) \left(\frac{t}{T_{\text{wrm}}} \right)$. Unless otherwise
 58 specified, we set $\eta_{\text{init}} = 0$ when referring to linear warmup.

59 **Parameterizations in Neural Networks:** The mechanism of warmup and its effectiveness is heavily
 60 influenced by the network parameterization (see Sections 3 and 4). Standard Parameterization (SP)
 61 [32] is a staple in common libraries [27, 3]. Another notable parameterization is the Neural Tangent
 62 Parameterization (NTP) [17], which along with SP resides in the kernel learning class at infinite
 63 width. Ref. [36] proposed Maximal Update Parameterization (μP) which exhibits feature learning at
 64 infinite width. Neural network parameterizations significantly impact training dynamics [19].

65 3 Warmup Mechanisms of Gradient and Adaptive Methods

66 This section analyzes the underlying mechanism of warmup through the lens of training instability.
 67 A key finding is a dichotomy between cooperative versus competitive dynamics based on how the
 68 natural evolution of the sharpness interplays with the training instability.

69 **3.1 Stochastic Gradient Descent (SGD)**

70 Learning rate warmup is intrinsically tied to sharpness dynamics, as sharpness determines the
 71 instability threshold η_c . As the learning rate is increased during warmup, these instabilities induce
 72 a temporary increase in the loss and a decrease in the sharpness to restore stability through the
 73 self-stabilization mechanism. Ultimately this allows the model to adapt to the increased learning
 74 rate. In other words, the primary goal of warmup is to gradually reduce sharpness, guiding training
 75 towards flatter regions that can accommodate training at higher learning rates [10].

76 However, digging deeper, we find that training has a ‘natural’ preference for sharpness evolution
 77 throughout the training course [19]. Before exceeding the instability threshold ($\eta < \eta_c$), training
 78 naturally experiences either a progressive increase or decrease in sharpness, as observed in Figure 1,
 79 which is unrelated to warmup. Here, the natural sharpness evolution can be defined as the change in
 80 sharpness experienced by gradient flow. The interplay between this natural sharpness evolution and
 81 the deliberate intervention of warmup to reduce sharpness can result in completely distinct dynamics.
 82 Below, we detail these cases and describe the conditions that typically exhibit them.

83 **(C1) Natural Progressive Sharpening** (top row of Figure 1): The combined effect of the network
 84 naturally increasing sharpness while the learning rate is also being increased results in a “head-on
 85 collision” at which the network reaches the instability threshold η_c . This causes the loss to increase,
 86 leading to a decrease in sharpness and facilitating a return to stability. As training proceeds, both
 87 sharpness and learning rate continue to increase, again surpassing the instability threshold. This
 88 results in a *persistent catapult cycle*, characterized by $\eta_t \approx 2/\lambda_t^H \approx \eta_c$, for the remainder of the
 89 warmup period, as seen in Figure 1(b).

90 **(C2) Natural Sharpness Reduction** (bottom row of Figure 1): The network is naturally already
 91 reducing its sharpness during early training. However, if the learning rate is increased sufficiently
 92 quickly, eventually the instability threshold will be reached (akin to a “rear-end collision”), causing
 93 the loss to increase. For small enough learning rates, the increased loss induces a dramatically more
 94 pronounced decrease in sharpness than would naturally occur, ultimately restoring stability. To
 95 exceed the instability threshold again, the learning rate must significantly increase to account for
 96 the decreased sharpness, potentially requiring considerable training steps. Consequently, training
 97 experiences one or more separated catapults during the warmup phase, as seen in Figure 1(c, d). This
 98 contrasts with the progressive sharpening case, where training enters a continuous catapult cycle after
 99 reaching the instability threshold for the first time. Notably, training may eventually reach a very
 100 flat region of the landscape during warmup, with gradients pointing towards increasing sharpness
 101 (e.g., $T_{\text{warm}} = 64$ in Figure 1(d)). Upon reaching such a region, the dynamics aligns with the natural
 102 progressive sharpening scenario.

103 **3.1.1 A Toy Model for Understanding the Warmup Mechanisms**

104 These two scenarios can be interpreted as cooperative or competitive dynamics between warmup
 105 and the natural evolution of sharpness. When training inherently undergoes sharpness reduction, it
 106 cooperates with warmup in decreasing sharpness. Conversely, if the natural trajectory of training is
 107 towards increasing sharpness, it opposes the warmup’s effort, leading to a persistent cycle of catapults.
 108 We can understand these mechanisms by analyzing a model of self-stabilization derived by [8].

109 The model assumes that the top eigenvector \mathbf{u} changes slowly through training and can be treated
 110 as constant. Next, consider a cubic approximation of the dynamics along a reference point θ^* . The
 111 dynamics along the projection $x_t := \mathbf{u}^T(\theta_t - \theta^*)$ is given by two coupled non-linear equations:

$$x_{t+1} = (1 - \eta_t \lambda_t^H) x_t, \quad \lambda_{t+1}^H = \lambda_t^H + \eta_t (\alpha - \beta x_t^2),$$

112 where $\alpha := -\nabla \lambda^H \cdot \nabla L(\theta)$ quantifies the instantaneous change in sharpness and $\beta := \|\nabla \lambda^H\|^2$
 113 controls to the non-linear change in sharpness. Ref. [8] considered a constant learning rate η
 114 and assumed progressive sharpening ($\alpha > 0$). Here, in contrast, we consider a time-dependent
 115 learning rate and allow α to attain both positive and negative values. In this model, an instability
 116 arises when $\eta_t \lambda_t^H > 2$. During instability, x_t continues to increase until the higher order term in
 117 the sharpness update equation causes a significant decrease in sharpness. Once the sharpness has
 118 decreased sufficiently, the stability is restored ($\eta_t \lambda_t^2 < 2$), and training continues.

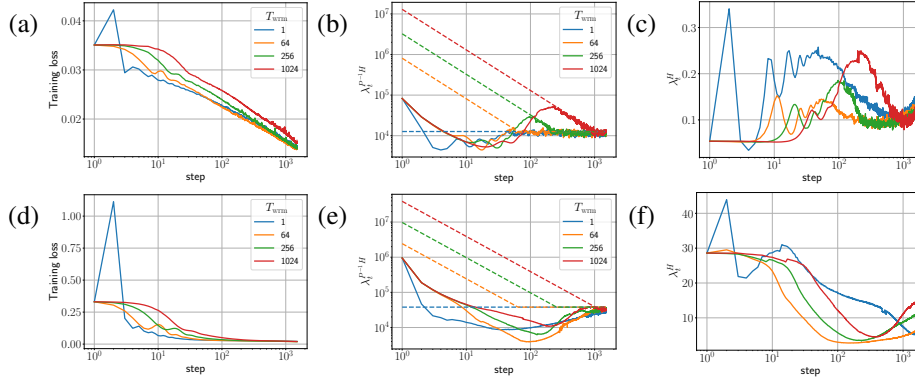


Figure 2: Training loss and sharpness trajectories of FCNs trained on the entire CIFAR-10 dataset with MSE loss using full batch Adam. (top) simple- μ P (for details, see Appendix D.2.1) with $\eta_{\text{trgt}} = 0.003$ and (bottom) SP with learning rate $\eta_{\text{trgt}} = 0.001$. The dashed lines in the sharpness figures illustrate the instability thresholds $(2+2\beta_1)/\eta_t(1-\beta_1)$. Similar mechanisms are observed for different architectures, loss functions, and smaller batch sizes as detailed in Appendix E.

119 Next, we consider the two natural sharpness evolution scenarios:

120 **(C1) Natural Progressive Sharpening ($\alpha > 0$):** The combined effect of naturally increasing
 121 sharpness ($\alpha > 0$) and the increasing learning rate from warmup leads to instability ($\eta_t \lambda_t^H > 2$).
 122 Resultantly, x_t increases until the higher order term in the sharpness update cause a decrease in
 123 sharpness ($x_t^2 > \frac{\alpha}{\beta}$). Once the sharpness has decreased appreciably so that $\eta_t \lambda_t^H < 2$, stability is
 124 restored and the training continues. As training proceeds, both progressive sharpening and increasing
 125 learning rate cause instability, resulting in a persistent catapult cycle characterized by $\eta_t \lambda_t^H \approx 2$.

126 **(C2) Natural Sharpness Reduction ($\alpha < 0$):** In this case, sharpness is naturally decreasing during
 127 training ($\alpha < 0$). If the learning rate is increased quick enough relative to decreasing sharpness, an
 128 instability occurs ($\eta_t \lambda_t^H > 2$). The increase in x_t causes a more pronounced decrease in sharpness
 129 than it would have occurred naturally, restoring instability. To exceed the instability threshold again,
 130 the learning rate must significantly increase to account for the decreased sharpness. This results in
 131 one or more separated catapults.

132 3.1.2 The Effect of Warmup Duration

133 Given a fixed target learning rate η_{trgt} , increasing the warmup duration T_{wrm} delays the point at which
 134 training exceeds the instability threshold η_c , allowing the sharpness to evolve freely before reaching
 135 this point. In the sharpness reduction case, sharpness can significantly decrease by the time this
 136 threshold is reached, lowering the need for warmup to decrease sharpness actively. Consequently,
 137 increasing T_{wrm} results in catapults that are both delayed and smaller in magnitude, as seen in
 138 Figure 1(d). As the catapults become less intense on increasing the warmup duration, the model
 139 can train at higher learning rates without diverging, pushing the divergence boundary. For extended
 140 warmup durations, warmup may not actively reduce sharpness in these sharpness reduction cases and
 141 instead “piggy-backs” on the inherent sharpness decrease.

142 In the progressive sharpening case, increasing T_{wrm} allows the sharpness to naturally increase. As a
 143 result, training exceeds the instability threshold for the first time at a relatively lower learning rate
 144 compared to the constant learning rate case. Although warmup has to now undertake more work in
 145 decreasing sharpness, it does so in a more gradual manner since increasing the warmup duration
 146 amounts to a lower warmup rate $\eta_{\text{trgt}}/T_{\text{wrm}}$. As a result, the fluctuations observed on exceeding the
 147 instability threshold are much smaller in magnitude, as seen in Figure 1(a, b).

148 3.1.3 Small vs. Large Initializations

149 So far, we have outlined different warmup mechanisms without describing specific conditions that
 150 typically exhibit them. Small initializations, such as those using maximal update parameterization
 151 (μ P) [36] or appropriately using normalizing layers (e.g. standard Transformer architectures, see
 152 Figure 14 in Appendix E.5), are characterized by a small initial network output. Such initializations
 153 start in flat regions where gradients point toward increasing sharpness [19], placing them in the

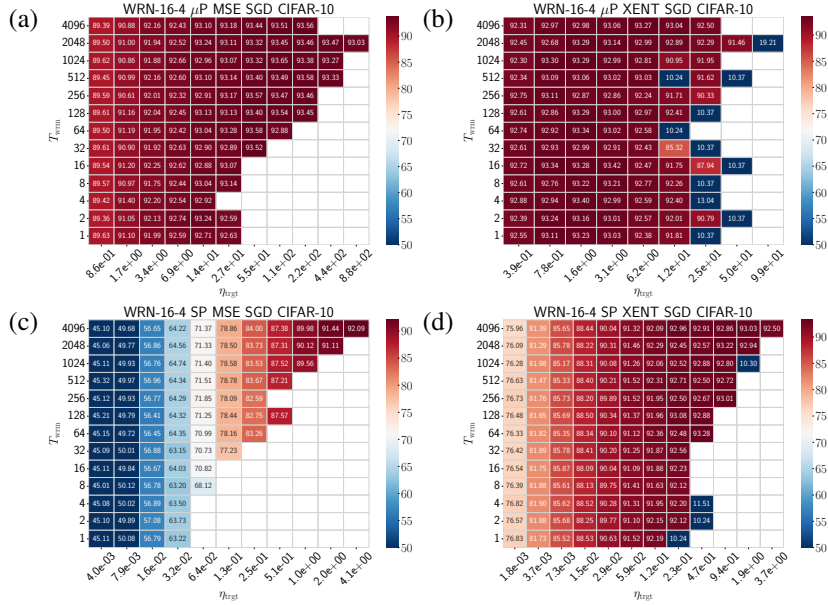


Figure 3: Test accuracy heatmaps of WideResNets (WRNs) trained on CIFAR-10 using different parameterizations and loss functions using SGD: (a) μ P and MSE loss, (b) μ P and cross-entropy loss, (c) SP and MSE loss, and (d) SP and cross-entropy loss. Empty cells correspond to training divergences. Additional results are provided in Appendix F.

154 progressive sharpening category (C1). As we will see in Section 4, such initializations may not
 155 significantly benefit from warmup as they already start in a flat region. In contrast, large initializations,
 156 such as FCNS, CNNs, ResNets with Standard Parameterization (SP) initialized at criticality [28, 30]
 157 or Transformers with the last layer-norm removed, undergo an early sharpness reduction, categorizing
 158 them into sharpness reduction category (C2). As the primary effect of warmup is to reduce sharpness,
 159 we expect such large initializations to considerably benefit from warmup. Notably, large initializations
 160 can eventually undergo progressive sharpening at later training stages [18, 19] and adhere to the
 161 second mechanism, especially for prolonged warmups. Instances of constant sharpness (C3) typically
 162 arise in models operating near the lazy regime [5], such as wide networks in NTP or SP.

163 **SGD with momentum:** The warmup mechanism of SGD with momentum, while at its core is
 164 similar to that of vanilla SGD, has a few subtleties. We discuss it in detail in Appendix E.2.

165 3.2 Adaptive Gradient Methods (Adam)

166 Figure 2 shows the training loss, pre-conditioned sharpness, and sharpness trajectories for full batch
 167 Adam. These results suggest that the local stability of adaptive optimizers is determined by the
 168 largest eigenvalue of the pre-conditioned Hessian, denoted by $\lambda^{P^{-1}H}$, rather than the sharpness itself
 169 (also, see Ref. [7] for late time instability). In these figures, sharpness is significantly smaller than
 170 its instability threshold $(2+2\beta_1)/\eta_t \approx 4000$, indicating that sharpness does not determine stability.
 171 Instead, loss catapults are associated with $\lambda^{P^{-1}H}$ exceeding its corresponding instability threshold.

172 The pre-conditioned sharpness starts high for both progressive sharpening (simple- μ P) and sharpness
 173 reduction (SP) scenarios considered in the previous section. For simplicity, we considered a simpler
 174 version of μ P, detailed in Appendix D.2.1. In particular, for μ P models, $\lambda_0^{P^{-1}H} \sim 10^5$ despite being
 175 initialized in a flat region as measured by sharpness, while for SP models, $\lambda_0^{P^{-1}H} \sim 10^6$. These large
 176 initial values of $\lambda_0^{P^{-1}H}$ can lead to training failures. We put forward strategies to improve Adam's
 177 initialization in Section 5; here we continue characterizing the warmup mechanisms of Adam.

178 Given that the pre-conditioned sharpness consistently starts high and decreases during early training,
 179 this behavior can be viewed as an extreme example of the natural sharpness reduction scenario (C2)
 180 described in the previous section. Training Adam at high initial learning rates without warmup can
 181 cause large catapults, as seen in Figure 2(d), potentially leading to training failures. Increasing the
 182 warmup duration allows the pre-conditioned sharpness to naturally decrease. This prevents the loss

183 from spiking during early training and avoids training failures. In the later stages of training, the
184 pre-conditioned sharpness may continue reducing or exhibit progressive sharpening. From here on,
185 the dynamics follows the warmup mechanisms discussed in the previous sections, with sharpness
186 replaced with pre-conditioned sharpness. Similar to the momentum case, Adam’s stability threshold
187 at late training times significantly decreases for smaller batch sizes [7], also shown in Appendix E.4.

188 4 Impact of Warmup on Training and Generalization

189 Here we investigate the impact of warmup on training and generalization by disentangling the role of
190 η_{tgt} and T_{wrm} . Our key findings are that generalization capability is primarily determined by η_{tgt} and
191 that Adam is particularly sensitive to large learning rates (specifically, large catapults). The role of
192 increasing T_{wrm} is to (i) allow the network to tolerate larger η_{tgt} , and (ii) move the network further
193 away from the divergence (failure) boundary, leading to a marginal improvement in generalization.
194 For experimental details, see Appendix D.

195 4.1 Stochastic Gradient Descent (SGD)

196 Figure 3 presents heatmaps that show the best test accuracy achieved during training, plotted in the
197 $\eta_{\text{tgt}}-T_{\text{wrm}}$ plane for different parameterizations and loss functions. These phase diagrams of warmup
198 also show the convergence-divergence boundary, with empty cells indicating training divergences,
199 illustrating the interplay between warmup duration and the maximum trainable η_{tgt} . Below, we
200 discuss the crucial insights these results provide into warmup’s role in training dynamics.

201 **Longer Warmup Facilitates Training at Higher Learning Rates:** These phase diagrams reveal
202 that an extended warmup duration facilitates training at higher target learning rates. This benefit is
203 particularly noticeable for large initializations (like SP) and MSE loss. In contrast, the advantage is
204 less pronounced when using cross-entropy loss and smaller initializations (like μP). The diminished
205 benefit for μP is likely due to its initialization in a relatively flat region of the loss landscape, which
206 can already facilitate training at higher learning rates at initialization. This consistent increase in
207 maximum η_{tgt} with warmup durations can be understood through the lens of warmup mechanisms
208 described in the previous section. As observed in Figure 1, when the warmup duration is increased,
209 loss catapults occurring on surpassing the instability thresholds become milder. This effectively
210 pushes the divergent boundary to higher learning rates.

211 **Final Performance Primarily Depends on the Target Learning Rate:** A closer look into these
212 phase diagrams reveals that, slightly away from the divergent boundary, the test accuracy primarily
213 depends on the target learning rate and nominally on the warmup duration. Based on the model
214 performance, we can categorize these phase diagrams into two distinct cases: (i) models that fail
215 to achieve optimal performance when trained with a constant learning rate (e.g., Figure 3(c)),
216 and (ii) models that attain optimal performance without warmup (e.g., Figure 3(b)). The first
217 scenario corresponds to models with large initializations. Increasing the warmup duration improves
218 performance by facilitating training at higher learning rates. Yet, similar performance is observed
219 for different warmup durations, suggesting that the primary gain comes from the target learning rate,
220 rather than the duration itself. The second case arises for flat initializations, which can already train
221 at large learning rates, and resultantly the optimal performance is already achieved without warmup.
222 While increasing warmup duration facilitates training at even higher learning rates, it does not enhance
223 performance. Nevertheless, it does broaden the range of optimal learning rates, reducing the need for
224 precise tuning of the target learning rate, and making training more practical and robust. We conclude
225 that warmup can serve two key purposes: (i) it can significantly improve model performance in large
226 initialization cases, and (ii) extend the range of optimal target learning rates for small initializations,
227 making it easier to tune the target learning rate. In Appendix F.2, we demonstrate that these results
228 hold on incorporating momentum and employing cosine learning rate decay.

229 4.2 Adam

230 The warmup phase diagrams for Adam, as shown in Figure 4(a), exhibit characteristics similar to the
231 sharpness reduction case of SGD, with notable differences. Increasing the warmup duration enables
232 training at higher learning rates by allowing the pre-conditioned sharpness to decrease naturally,
233 thereby reducing the severity of catapults. These large catapults, which may persist in Adam’s

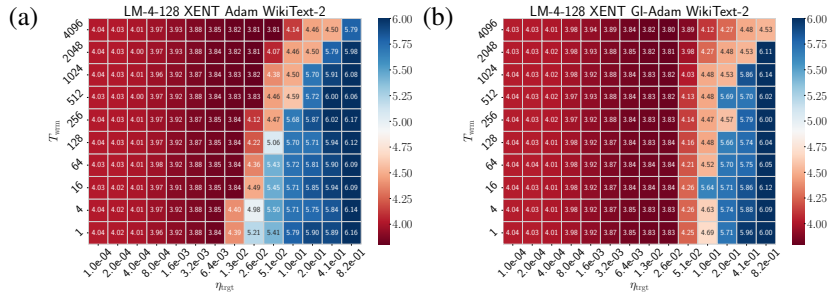


Figure 4: Test loss heatmaps of Pre-LN Transformers in SP trained on WikiText-2 with cross-entropy loss using (a) Adam, and (b) GI-Adam (introduced in Section 5).

memory, can lead to performance degradation and training failures. Thus, in addition to facilitating training at higher rates similar to SGD, warmup further improves Adam’s performance by addressing its vulnerability to large catapults, justifying its widespread use with Adam. Below, we discuss the distinct properties of Adam phase diagrams in detail.

Training Failures of Adam: Remarkably, we find that models trained with Adam always exhibit training failures rather than divergences where the loss grows without bound, as further demonstrated in Appendix G. In cases of training failure, we often observed that certain layers or residual blocks output zero, leading to vanishing gradients. This implies that the model gets stuck at a critical point and is unable to train further. Understanding this unexpected phenomenon requires further study, which we leave to future work.

Performance Degradation prior to Failure Boundary: Test accuracy in these phase diagrams declines well before the failure boundary, in stark contrast to SGD where optimal learning rates are observed near the divergence boundary. This discrepancy stems from Adam’s property of retaining a memory of gradient magnitudes. At large learning rates, along with the loss, the gradients spike during early training, as seen in Figure 23 in Appendix G. While the gradients decrease after a few training steps, the second moment of gradients v remains large for an extended period, leading to a small effective learning rate ηP^{-1} . As a result, training struggles to escape high-loss regions. Therefore, a longer warmup is more beneficial for Adam compared to SGD, as it is crucial to stay away from the failure boundary.

5 GI-Adam: Improving Adam’s Initialization

In Section 3.2, we observed that the pre-conditioned sharpness for Adam starts at a high value, even for low sharpness initializations like μP , and can lead to training failures at large learning rates. We propose Gradient Initialized Adam (GI-Adam), which initializes the second moment using the gradient squared, $v_0 = g_0^2$. In Appendix H.2, we show that a bias correction is not required when the second moment is initialized using the gradients. As a result, GI-Adam can be viewed as standard Adam with an automated warmup given by $\eta_t = \eta_{\text{tgt}} \sqrt{1 - \beta_2^t}$.

This simple trick reduces the initial pre-conditioned sharpness by around two orders of magnitude (more precisely by a factor of $\sqrt{1 - \beta_2}$) at initialization, preventing large catapults, as illustrated in Figure 25 of Appendix H.1 (c.f. Figure 2(d-f)). Moreover, it consistently shows improvement over standard Adam across datasets and prevents training failures by pushing the training failure boundary to higher η_{tgt} , as shown in Figure 4(b). We provide additional results for different datasets in Appendix F.3. To further assess that the primary cause of instability during early training is the large pre-conditioned sharpness, we randomly initialize v_0 but with the same norm as the gradients at initialization. Like GI-Adam, this also results in improved performance as shown in Appendix H.3.

6 Discussion

Our analysis provides new insights into the role of warmup across optimizers and parameterizations. We found compelling evidence that the primary effect of warmup is to facilitate training at higher learning rates and stabilizing the training dynamics by keeping it away from the failure (divergence) boundary. Looking under the hood, we found a variety of underlying mechanisms, which also suggested several improvements for hyperparameter initialization. In Appendix A we provide practical guidance for practitioners on choosing the warmup duration.

275 References

- 276 [1] Tiny imagenet challenge. <https://cs231n.stanford.edu/reports/2017/pdfs/930.pdf>.
277 pdf.
- 278 [2] Stephen P Boyd and Lieven Vandenbergh. *Convex optimization*. Cambridge university press,
279 2004.
- 280 [3] James Bradbury, Roy Frostig, Peter Hawkins, Matthew James Johnson, Chris Leary, Dougal
281 Maclaurin, George Necula, Adam Paszke, Jake VanderPlas, Skye Wanderman-Milne, and Qiao
282 Zhang. JAX: composable transformations of Python+NumPy programs, 2018.
- 283 [4] Sébastien Bubeck et al. Convex optimization: Algorithms and complexity. *Foundations and
284 Trends® in Machine Learning*, 8(3-4):231–357, 2015.
- 285 [5] Lénaïc Chizat, Edouard Oyallon, and Francis Bach. On lazy training in differentiable program-
286 ming. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett,
287 editors, *Advances in Neural Information Processing Systems*, volume 32. Curran Associates,
288 Inc., 2019.
- 289 [6] Jeremy Cohen, Simran Kaur, Yuanzhi Li, J Zico Kolter, and Ameet Talwalkar. Gradient descent
290 on neural networks typically occurs at the edge of stability. In *International Conference on
291 Learning Representations*, 2021.
- 292 [7] Jeremy M. Cohen, Behrooz Ghorbani, Shankar Krishnan, Naman Agarwal, Sourabh Medapati,
293 Michal Badura, Daniel Suo, David Cardoze, Zachary Nado, George E. Dahl, and Justin Gilmer.
294 Adaptive gradient methods at the edge of stability, 2022.
- 295 [8] Alex Damian, Eshaan Nichani, and Jason D. Lee. Self-stabilization: The implicit bias of
296 gradient descent at the edge of stability. In *The Eleventh International Conference on Learning
297 Representations*, 2023.
- 298 [9] Emily Dinan, Sho Yaida, and Susan Zhang. Effective theory of transformers at initialization,
299 2023.
- 300 [10] Justin Gilmer, Behrooz Ghorbani, Ankush Garg, Sneha Kudugunta, Behnam Neyshabur, David
301 Cardoze, George Edward Dahl, Zachary Nado, and Orhan Firat. A loss curvature perspective
302 on training instabilities of deep learning models. In *International Conference on Learning
303 Representations*, 2022.
- 304 [11] Gabriel Goh. Why momentum really works. *Distill*, 2017.
- 305 [12] Akhilesh Gotmare, Nitish Shirish Keskar, Caiming Xiong, and Richard Socher. A closer look
306 at deep learning heuristics: Learning rate restarts, warmup and distillation. In *International
307 Conference on Learning Representations*, 2019.
- 308 [13] Priya Goyal, Piotr Dollár, Ross Girshick, Pieter Noordhuis, Lukasz Wesolowski, Aapo Kyrola,
309 Andrew Tulloch, Yangqing Jia, and Kaiming He. Accurate, large minibatch sgd: Training
310 imagenet in 1 hour. *arXiv preprint arXiv:1706.02677*, 2017.
- 311 [14] Kaiming He, X. Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image
312 recognition. *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*,
313 pages 770–778, 2016.
- 314 [15] Jonathan Heek, Anselm Levskaya, Avital Oliver, Marvin Ritter, Bertrand Rondepierre, Andreas
315 Steiner, and Marc van Zee. Flax: A neural network library and ecosystem for JAX, 2020.
- 316 [16] Dan Hendrycks and Kevin Gimpel. Bridging nonlinearities and stochastic regularizers with
317 gaussian error linear units. *CoRR*, abs/1606.08415, 2016.
- 318 [17] Arthur Jacot, Franck Gabriel, and Clement Hongler. Neural tangent kernel: Convergence and
319 generalization in neural networks. In S. Bengio, H. Wallach, H. Larochelle, K. Grauman,
320 N. Cesa-Bianchi, and R. Garnett, editors, *Advances in Neural Information Processing Systems*,
321 volume 31. Curran Associates, Inc., 2018.

- 322 [18] Dayal Singh Kalra and Maissam Barkeshli. Phase diagram of early training dynamics in deep
323 neural networks: effect of the learning rate, depth, and width. In *Thirty-seventh Conference on*
324 *Neural Information Processing Systems*, 2023.
- 325 [19] Dayal Singh Kalra, Tianyu He, and Maissam Barkeshli. Universal sharpness dynamics in neural
326 network training: Fixed point analysis, edge of stability, and route to chaos. *arXiv preprint*
327 *arXiv:2311.02076*, 2023.
- 328 [20] Diederik Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *International*
329 *Conference on Learning Representations (ICLR)*, San Diego, CA, USA, 2015.
- 330 [21] Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images.
331 2009.
- 332 [22] Aitor Lewkowycz, Yasaman Bahri, Ethan Dyer, Jascha Sohl-Dickstein, and Guy Gur-Ari.
333 The large learning rate phase of deep learning: the catapult mechanism. *arXiv preprint*
334 *arXiv:2003.02218*, 2020.
- 335 [23] Liyuan Liu, Haoming Jiang, Pengcheng He, Weizhu Chen, Xiaodong Liu, Jianfeng Gao, and
336 Jiawei Han. On the variance of the adaptive learning rate and beyond. In *International*
337 *Conference on Learning Representations*, 2020.
- 338 [24] Ilya Loshchilov and Frank Hutter. SGDR: Stochastic gradient descent with warm restarts. In
339 *International Conference on Learning Representations*, 2017.
- 340 [25] Stephen Merity, Caiming Xiong, James Bradbury, and Richard Socher. Pointer sentinel mixture
341 models. In *International Conference on Learning Representations*, 2017.
- 342 [26] Kevin P Murphy. *Probabilistic machine learning: an introduction*. MIT press, 2022.
- 343 [27] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan,
344 Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. Pytorch: An imperative
345 style, high-performance deep learning library. *Advances in neural information processing*
346 *systems*, 32, 2019.
- 347 [28] Ben Poole, Subhaneil Lahiri, Maithra Raghu, Jascha Sohl-Dickstein, and Surya Ganguli.
348 Exponential expressivity in deep neural networks through transient chaos. In *NIPS*, 2016.
- 349 [29] Alec Radford, Jeff Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. Language
350 models are unsupervised multitask learners. 2019.
- 351 [30] Daniel A. Roberts, Sho Yaida, and Boris Hanin. *The Principles of Deep Learning Theory*.
352 Cambridge University Press, 2022. <https://deeplearningtheory.com>.
- 353 [31] Rico Sennrich, Barry Haddow, and Alexandra Birch. Neural machine translation of rare words
354 with subword units. In Katrin Erk and Noah A. Smith, editors, *Proceedings of the 54th Annual*
355 *Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages
356 1715–1725, Berlin, Germany, August 2016. Association for Computational Linguistics.
- 357 [32] Jascha Narain Sohl-Dickstein, Roman Novak, Samuel S. Schoenholz, and Jaehoon Lee. On the
358 infinite width limit of neural networks with a standard parameterization. *ArXiv*, abs/2001.07301,
359 2020.
- 360 [33] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez,
361 Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information*
362 *processing systems*, 30, 2017.
- 363 [34] Lei Wu, Chao Ma, et al. How sgd selects the global minima in over-parameterized learning: A
364 dynamical stability perspective. *Advances in Neural Information Processing Systems*, 31, 2018.
- 365 [35] Ruibin Xiong, Yunchang Yang, Di He, Kai Zheng, Shuxin Zheng, Chen Xing, Huishuai
366 Zhang, Yanyan Lan, Liwei Wang, and Tieyan Liu. On layer normalization in the transformer
367 architecture. In Hal Daumé III and Aarti Singh, editors, *Proceedings of the 37th International*
368 *Conference on Machine Learning*, volume 119 of *Proceedings of Machine Learning Research*,
369 pages 10524–10533. PMLR, 13–18 Jul 2020.

- 370 [36] Greg Yang and Edward J. Hu. Tensor programs iv: Feature learning in infinite-width neural
 371 networks. In Marina Meila and Tong Zhang, editors, *Proceedings of the 38th International
 372 Conference on Machine Learning*, volume 139 of *Proceedings of Machine Learning Research*,
 373 pages 11727–11737. PMLR, 18–24 Jul 2021.
- 374 [37] Greg Yang, Edward J Hu, Igor Babuschkin, Szymon Sidor, Xiaodong Liu, David Farhi, Nick
 375 Ryder, Jakub Pachocki, Weizhu Chen, and Jianfeng Gao. Tuning large neural networks via
 376 zero-shot hyperparameter transfer. In A. Beygelzimer, Y. Dauphin, P. Liang, and J. Wortman
 377 Vaughan, editors, *Advances in Neural Information Processing Systems*, 2021.
- 378 [38] Sergey Zagoruyko and Nikos Komodakis. Wide residual networks. 2017.
- 379 [39] Aston Zhang, Zachary C Lipton, Mu Li, and Alexander J Smola. *Dive into deep learning*.
 380 Cambridge University Press, 2023.
- 381 [40] Hongyi Zhang, Moustapha Cisse, Yann N. Dauphin, and David Lopez-Paz. mixup: Beyond
 382 empirical risk minimization. In *International Conference on Learning Representations*, 2018.

383 A Practical Guidance for Practitioners

384 **How to Select the Warmup Duration?** Given a target learning rate η_{tgt} , if the training loss during
 385 the warmup period exhibits large instabilities (loss spikes), the warmup duration T_{wrm} should be
 386 increased until such instabilities are sufficiently small. This effectively moves training away from the
 387 divergent / failure boundary, as illustrated in Figure 3. This is particularly crucial for Adam, as large
 388 instabilities can be detrimental and lead to considerable performance degradation without divergence,
 389 as discussed in Section 4.2.

390 **How to Select the Target Learning Rate?** As the primary effect of warmup is to anneal sharpness
 391 by increasing the learning rate beyond the instability threshold, it suggests that the target learning
 392 rate should be at least greater than the instability threshold at initialization.

393 **When to Decay the Learning Rate?** Figure 18 suggests that employing learning rate decay at
 394 small learning rates can result in performance degradation for a fixed training budget. Therefore, the
 395 learning rate should be decayed at large target learning rates only. The underlying intuition is that we
 396 use large target learning rates to train in a flat region of the landscape. However, these large learning
 397 rates restrict training to go into sharper regions of the basin and learning rate decay helps.

398 **Leveraging μP for Effecient Training:** Our analysis suggests that the primary role of warmup
 399 facilitates training at higher learning rates by gradually reducing sharpness. Given this perspective,
 400 beginning training with flat initializations, such as μP , is advantageous. These initializations might
 401 allow for achieving optimal performance without the need for warmup, as observed in Figure 3.

402 B Overview of Training Instabilities and the Self-Stabilization Mechanism

403 The underlying mechanism of warmup is intimately tied to training instabilities. These training insta-
 404 bilities, often referred to as ‘catapults’ [22, 6], arise when the learning rate η exceeds a critical thresh-
 405 old η_c , where both η and η_c generally change with time. When the instability threshold is exceeded
 406 ($\eta > \eta_c$), two cases arise: (i) if the learning rate is higher than the instability threshold but smaller
 407 than a maximum stable learning rate (which varies with time), i.e., $\eta_c < \eta < \eta_{\text{max}}$, training stabilizes
 408 through a self-stabilization process and training continues, (ii) if the learning rate exceeds this maxi-
 409 mum stable learning rate $\eta > \eta_{\text{max}}$, training experiences severe instabilities. For SGD, these can result
 410 in training divergence, characterized by the loss increasing to infinity, whereas for Adam, training may
 411 cease, resulting in a training failure, where the loss fails to improve significantly over its initial value.

412 For vanilla GD, the critical threshold is related to sharpness as $\eta_c \approx 2/\lambda^H$ ¹, and the self-stabilization
 413 mechanism can be described as a four-step process [22, 8]. To illustrate this, consider the $T_{\text{wrm}} = 64$

¹This relationship holds for the MSE loss and simple settings only. For an overview of instability thresholds in various settings and different optimizers, see Appendix C.1.

trajectories depicted in Figure 1(c, d). In the sharpness plot, the dashed lines represent the $2/\eta_t$ curves, and when λ_t^H is above these curves, training exceeds the instability threshold ($\eta > \eta_c$). The four steps of the self-stabilization mechanism are:

- (1) **Approaching instability:** Due to increasing learning rate and/or progressive sharpening, training approaches the instability threshold $\eta = \eta_c$. In Figure 1(d), this occurs within the first 10 steps due to increasing learning rate.
- (2) **Loss increases:** The loss begins to rise when the instability threshold is exceeded ($\eta > \eta_c$), as seen in Figure 1(c).
- (3) **Sharpness reduction:** For small enough learning rates, the increasing loss causes an abrupt decrease in sharpness, as observed in Figure 1(d). If the sharpness fails to decrease over extended steps, it may result in training divergence (e.g., see $T_{\text{wrm}} = 1$ trajectories in the same figure).
- (4) **Return to stability:** The reduction in sharpness causes $\eta_c = 2/\lambda^H$ to increase, restoring stability ($\eta < \eta_c$) and allowing for an eventual loss decrease.

While the self-stabilization process for more complex optimizers, such as SGD with momentum or Adam, remains poorly understood, a qualitatively similar mechanism is observed in practice, as we will see in the later sections.

The critical learning rate η_c is influenced by a variety of factors, including the choice optimizer [6, 7], mini-batch size [34, 7], and model properties such as depth, width, parameterization, and initialization [18, 19]. For a detailed overview of instability thresholds, see Appendix C.

C Instability Thresholds

C.1 Overview of Instability Thresholds

Lewkowycz et al. [22] showed that for wide networks in NTP/SP trained with MSE loss and SGD, this critical learning rate is $2/\lambda_0^H$ early in training. Further investigation by Kalra and Barkeshli [18] demonstrated that sharpness reduction during early training causes η_c to increase with depth and $1/\text{width}$. In such scenarios, η_c can be as large as $40/\lambda_0^H$. Cohen et al. [6] demonstrated that sharpness at late training times for GD with momentum coefficient β oscillates above $(2+2\beta)/\eta$, suggesting $\eta_c \gtrsim (2+2\beta)/\lambda_t^H$ at late training times. Expanding on this, Cohen et al. [7] analyzed adaptive optimizers and found that for Adam, the pre-conditioned sharpness $\lambda^{P-1}H$ oscillates around $(2+2\beta_1)/\eta(1-\beta_1)$ at late training times. The instability threshold also depends on the mini-batch size [34] and is often observed to be smaller than their full batch counterparts [6, 7].

D Experimental Details

This section provides additional experimental details. All models were implemented using the JAX [3], and Flax libraries [15]. The key results can be reproduced using the GitHub repo: <https://github.com/dayal-kalra/why-warmup>.

Experimental Setup for Section 4: We consider WideResNets (WRNs) and Transformers (LM) parameterized in either SP or μ P. WRNs are trained on standard classification tasks such as CIFAR-10, CIFAR-100, and Tiny-ImageNet, employing data augmentation. Transformers are trained on the next token prediction task using the WikiText-2 dataset. These models are trained with MSE or cross-entropy (xent) loss functions using SGD or Adam optimizers for a fixed training budget of $T = 10^5$ steps unless otherwise specified. Training begins with a linear warmup phase from $\eta_{\text{init}} = 0$ to η_{trgt} over T_{wrm} steps. After the warmup phase, training continues at η_{trgt} for the remaining training budget. In some cases, following the warmup period, we gradually decrease the learning rate using cosine decay [24]. Target learning rates are sampled exponentially until divergence or a ‘training failure’ is observed. Here, training failure refers to instances where the performance at the end of the training fails to improve significantly compared to its initial value. For example, if the final training accuracy for a classification task is less than 1.5 times the accuracy of a random guess, we consider it as a training failure. We refer to the transition between convergence and training failure as the failure boundary. Further details are provided in Appendix D.

463 D.1 Datasets Details

464 D.1.1 Image Classification Tasks

465 We consider standard image classification datasets such as CIFAR-10, CIFAR-100 [21], and Tiny-
466 ImageNet [1]. The images are normalized to have zero mean and unit variance. For MSE loss, we
467 use one-hot encoding for the labels.

468 **Data augmentation:** For various image classification tasks, we employ data augmentation techniques,
469 applied in the following order: random horizontal flips, random cropping, and mixup [40].

470 D.1.2 Language Modeling Tasks

471 We consider the next token prediction task on the Wikitext-2 dataset [25], consisting of $\sim 2M$ tokens.
472 We use Byte Pair Encoding (BPE) tokenizer [31] with a Whitespace pre-tokenizer. Due to the high
473 computational cost associated with hyperparameter tuning, we restrict to smaller models with $\sim 2M$
474 parameters. Furthermore, we restrict the vocabulary size to 4096 to ensure that embedding parameters
475 do not dominate the total number of parameters in the model.

476 D.2 Model Details

477 This section describes the models considered, including their parameterization and initialization
478 details. We adopt parameterizations outlined in Table 9 of Ref. [37]. Unless otherwise specified, we
479 employ ReLU non-linearities and initialize the weights with a truncated normal distribution², with a
480 variance $\sigma_w^2 = 2.0$ in appropriate parameterizations (details below), except for the last layer, which
481 has a weight variance of $\sigma_w^2 = 1.0$. All biases are initialized to zeros.

482 D.2.1 Parameterizations

483 **Standard Parameterization (SP):** For SP, the weights are initialized with truncated Gaussian
484 distribution $\mathcal{N}(0, \sigma_w^2/f_{\text{an}_{\text{in}}})$ and the biases are initialized to zero.

485 **Maximal Update Parameterization (μP):** For μP , different schemes are employed for the inter-
486 mediate and last layers. The intermediate layers are initialized using $\mathcal{N}(0, \sigma_w^2/f_{\text{an}_{\text{out}}})$ and the layer
487 outputs are scaled by the factor $\sqrt{f_{\text{an}_{\text{out}}}/f_{\text{an}_{\text{in}}}}$. In comparison, the layer weights are initialized with
488 $\mathcal{N}(0, \sigma_w^2/f_{\text{an}_{\text{in}}})$, and the final output is rescaled by the factor $\sqrt{1/f_{\text{an}_{\text{in}}}}$. Conveniently, for SGD, the
489 learning rate does not scale with width in the above μP formulation. In comparison, for Adam,
490 the learning rate corresponding to input, intermediate, and output layers are rescaled by the factors
491 $1/\sqrt{f_{\text{an}_{\text{out}}}}$, $1/\sqrt{f_{\text{an}_{\text{in}}}}$ and $1/f_{\text{an}_{\text{in}}}$. Since we are utilizing μP only to obtain flat initializations, we omit the
492 additional scaling of the learning rate for Adam in some experiments (e.g., Figure 2). As a result, the
493 instability threshold is only dependent on the target learning rate η_{trgt} during late training, rather than
494 on the largest learning rate across layers. We refer to this parameterization as ‘simple- μP ’ for Adam.

495 D.2.2 Architectures

496 **Fully Connected Networks (FCNs):** We consider fully connected networks with a constant width
497 of n and a depth of d layers. These networks are denoted by FCN- d - n . Unless specified, we
498 considered $d = 4$ layer FCNs with width $n = 512$.

499 **WideResNets (WRNs):** We consider WideResNets [38] with d layers, S stages, and a widening
500 factor of k , denoted by WRN- d - k . The number of channels in each stage $s \in [0, S)$ is given
501 by $2^s \times 16 \times k$, with the input layer having 16 channels. For example, WRN-16-4 consists of
502 $S = 3$ stages, each with $[2, 2, 2]$ layers, and the corresponding number of channels in each stage is
503 $[64, 128, 256]$. In all our experiments, we use LayerNorm instead of BatchNorm.

504 **Transformers:** We consider Transformers with GPT-2 style architecture [29]. These models use
505 sinusoidal positional embeddings [33] and are implemented in the Standard Parameterization (SP)

²for details, see https://jax.readthedocs.io/en/latest/_autosummary/jax.nn.initializers.truncated_normal.html

506 with GELU activation [16]. We initialize all layers using the $\sigma_w^2/\text{fan}_{\text{in}}$ scheme, except for the embedding
 507 layers, as they do not involve matrix multiplication [9]. We consider both Pre-LN [35] and Post-LN
 508 [33] Transformer variants. We denote a Transformer with d blocks and an embedding dimension of n
 509 as LM- d - n . Unless specified, the model has $d = 4$ blocks, embedding dimension $n = 128$, context
 510 length $T_{\text{ctxt}} = 64$ and are trained for 10^4 steps.

511 D.3 Optimization Details

512 D.3.1 Optimizers

513 **SGD(-M):** Given gradients \mathbf{g}_t at step t , Stochastic Gradient Descent with momentum (SGD-M)
 514 updates the parameters $\boldsymbol{\theta}_t$ using learning rate η_t and momentum \mathbf{m}_t with coefficient β . The update
 515 equations are:

$$\mathbf{m}_t = \mathbf{g}_t + \beta \mathbf{m}_{t-1}, \quad (1)$$

$$\boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t - \eta_t \mathbf{m}_t. \quad (2)$$

516 Here, $\beta = 0$ corresponds to SGD. In all experiments incorporating momentum, the default value of
 517 the coefficient is set to $\beta = 0.9$.

518 **Adam:** Given gradients \mathbf{g}_t at step t , Adam [20] updates the parameters $\boldsymbol{\theta}_t$ using learning rate η_t
 519 and the first two moments of the gradient \mathbf{m}_t and \mathbf{v}_t with their coefficients β_1 and β_2 , respectively.
 520 The equations governing the updates are:

$$\mathbf{m}_t = \beta_1 \mathbf{m}_{t-1} + (1 - \beta_1) \mathbf{g}_t, \quad (3)$$

$$\mathbf{v}_t = \beta_2 \mathbf{v}_{t-1} + (1 - \beta_2) \mathbf{g}_t^2, \quad (4)$$

$$\boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t - \eta_t \frac{\hat{\mathbf{m}}_t}{\sqrt{\hat{\mathbf{v}}_t + \epsilon}}, \quad (5)$$

521 where $\hat{\mathbf{m}}_t = \frac{\mathbf{m}_t}{1 - \beta_1^t}$ and $\hat{\mathbf{v}}_t = \frac{\mathbf{v}_t}{1 - \beta_2^t}$ are the bias-corrected moments, and ϵ is a small scalar used for
 522 numerical stability. The pre-conditioner for Adam is given by:

$$P_t = (1 - \beta_1^t) \left[\text{diag} \left(\frac{\mathbf{v}_t}{1 - \beta_2^t} \right) + \epsilon \mathbf{I} \right]. \quad (6)$$

523 In all experiments, the default values are set to $\beta_1 = 0.9$, $\beta_2 = 0.999$, and $\epsilon = 10^{-8}$, unless otherwise
 524 specified.

525 D.3.2 Linear Warmup

526 Warmup linearly increases the learning rate from an initial value η_{init} to a target value η_{trgt} over T_{wrm}
 527 training steps. The learning rate η_t at step t is given by:

$$\eta_t = \eta_{\text{init}} + (\eta_{\text{trgt}} - \eta_{\text{init}}) \left(\frac{t}{T_{\text{wrm}}} \right). \quad (7)$$

528 Here, $\alpha := \frac{(\eta_{\text{trgt}} - \eta_{\text{init}})}{T_{\text{wrm}}}$ is referred to as the rate of warmup. Under the above definition, constant
 529 learning rate training corresponds to $T_{\text{wrm}} = 1$. $T_{\text{wrm}} = 1$ corresponds to constant learning rate.
 530 Unless otherwise specified, we set $\eta_{\text{init}} = 0$ when referring to linear warmup.

531 D.3.3 Learning Rate Decay

532 In several experiments, we employ learning rate decay following the warmup phase. Specifically, we
 533 use cosine learning rate decay, which is detailed below.

534 **Cosine Decay:** Towards the end of training, it is typical to reduce the learning rate to a small value.
 535 Cosine decay is a commonly used method for decaying the learning rate from an initial value of η_{trgt}
 536 down to a value η_{min} over T_{cos} steps, according to the rule:

$$\eta_t = \eta_{\text{trgt}} + (\eta_{\text{min}} - \eta_{\text{trgt}}) \left[\frac{1}{2} \left(1 + \cos \left(\frac{\pi t}{T_{\text{cos}}} \right) \right) \right]^\rho, \quad (8)$$

537 where ρ governs the rate of decay, with $\rho = 1$ being the standard. Note that with $\rho = 0$, the learning
 538 rate is not decayed and instead maintained at η_{trgt} . In the above expression, t counts the steps from
 539 the initiation of cosine decay and not the current training step. As per standard practice, we consider
 540 $\rho = 1$ and decay the learning rate to $\eta_{\text{min}} = \eta_{\text{trgt}}/10$.

541 D.3.4 Target Learning Rate Sampling for Phase Diagrams

542 For SGD, target learning rates η_{trgt} are exponentially sampled using the initial sharpness λ_0^H . Starting
 543 with $\eta_{\text{trgt}} = 1/\lambda_0^H$, subsequent rates are sampled until divergence as $2^x/\lambda_0^H$ for values of x increased
 544 in integer steps starting from zero. For WRNs trained with Adam, we sample target learning rates
 545 exponentially as $\eta_{\text{trgt}} = 2^x \times 10^{-5}$, where x is incremented in integer steps starting from zero until
 546 training failure. For Transformers, we sample the learning rate in a similar fashion but starting from
 547 10^{-4} and increment x in steps of 0.5.

548 D.4 Sharpness and Pre-conditioned Sharpness Measurement

549 We measured sharpness / pre-conditioned sharpness using the JAX implementation of the LOBPCG
 550 sparse eigenvalue solver with the tolerance set to 10^{-9} and maximum number of iterations to $n_{\text{iter}} =$
 551 1000. In most cases, the solver converges within 40 iterations. We performed these computations in
 552 float64, as the solver would not converge with float32 in some cases.

553 In certain instances, the pre-conditioned sharpness computation did not converge within 1000 solver
 554 iterations. Moreover, we observed that the solver converges on restarting it with a new initial guess
 555 of the eigenvector within 40 iterations. To address these edge cases, we employed the following
 556 method: if the solver did not converge within 100 iterations, we restarted it with a new initial guess
 557 for the eigenvector. We allowed for at most 10 restarts with the maximum number of iterations set to
 558 $n_{\text{iter}} = 1000$ in the last attempt. In all reported cases, the solver converges using this method.

559 D.5 Additional Figure Details

560 **Figure 1:** Training trajectories of 4-layer FCNs with width $n = 512$, trained on a 5k subset of
 561 CIFAR-10 using MSE loss and GD in (top) μP with $\eta_{\text{trgt}} = 1/\lambda_0^H$, where $\lambda_0^H \approx 0.05$, and (bottom) SP
 562 with $\eta_{\text{trgt}} = 32/\lambda_0^H$, where $\lambda_0^H \approx 50$.

563 **Figure 2:** Training loss and sharpness trajectories of 4 layer FCNs with width $n = 512$, in (top) μP
 564 with learning rate $\eta_{\text{trgt}} = 0.003$ and (bottom) SP with $\eta_{\text{trgt}} = 0.001$ trained the CIFAR-10 dataset with
 565 MSE loss using full batch Adam with $\beta_1 = 0.9$, $\beta_2 = 0.999$ and $\epsilon = 10^{-8}$. In these experiments, we
 566 use data augmentation as described in Appendix D.1.1.

567 **Figure 3:** Test accuracy heatmaps of WRN-16-4 trained on CIFAR-10 using different parameter-
 568 izations and loss functions using SGD with a batch size $B = 128$: (a) SP and MSE loss, (b) μP
 569 and cross-entropy loss (c) SP and cross-entropy loss. All models are trained for 10^5 steps. In these
 570 experiments, we use data augmentation as described in Appendix D.1.1.

571 **Figure 4:** Test loss heatmaps of Pre-LN Transformers in SP trained on WikiText-2 with cross-
 572 entropy loss using (a) Adam, and (b) GI-Adam (introduced in Section 5) over Adam. The Transformer
 573 models have $d = 4$ blocks, embedding dimension $n = 128$, a context length of $T_{\text{cnxt}} = 64$. These
 574 experiments also employ cosine decay, as described in Appendix D.3.3.

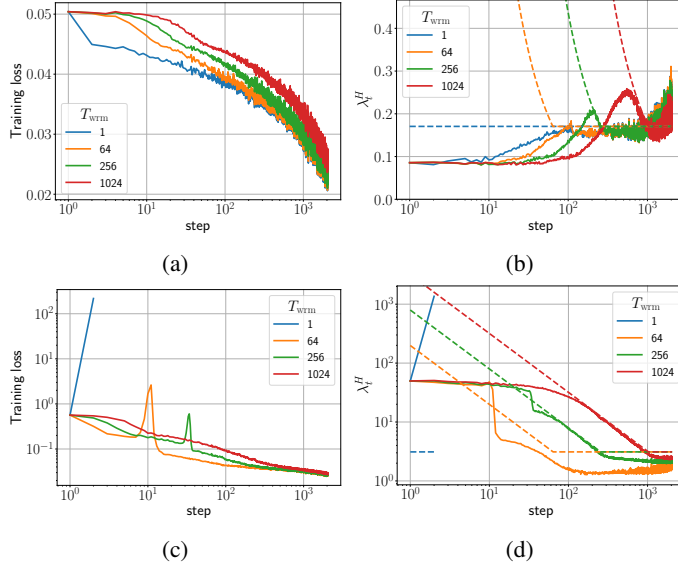


Figure 5: Training loss and sharpness trajectories of FCNs trained on CIFAR-10 with MSE loss using SGD with a batch size $B = 512$. The dashed lines in the sharpness figures illustrate the instability thresholds $2/\eta_t$. (top) μP with learning rate $1/\lambda_0^H$, (bottom) SP with learning rate $32/\lambda_0^H$.

575 D.6 Estimation of Computational Resources

576 The phase diagram experiments typically required about an hour on per run on an A100 GPU.
 577 Consequently, each phase diagram consumed approximately 100 A100 hours of computational time.
 578 With a total of 16 phase diagrams, this equates to 1600 A100 hours dedicated solely to phase diagram
 579 computations. Additionally, the warmup mechanism experiments, which were conducted over 2000
 580 steps, required sharpness estimation. The FCN experiments required approximately 1200 A100 hours,
 581 while the WRN mechanism experiments consumed 1600 A100 hours. The experiments concerning
 582 the initial learning rate took about 20 A100 hours. This brings the total computational time amounted
 583 to approximately 4500 A100 hours. Preliminary experiments took about 1000 A100 hours. Hence,
 584 we estimate the total computational cost to be around 5500 A100 hours.

585 E Additional Results for Mechanisms of Warmup

586 This section presents additional trajectories for warmup mechanisms discussed in Section 3 covering
 587 various architectures, loss functions, and optimizers.

588 E.1 Stochastic Gradient Descent

589 Figure 5 shows that the warmup mechanisms for full batch GD are also observed in the SGD with a
 590 batch size $B = 512$. The results for other optimizers in the mini-batch setting are discussed in their
 591 respective sections.

592 E.2 Stochastic Gradient Descent with Momentum

593 While the warmup mechanisms of SGD with momentum are fundamentally similar to those of vanilla
 594 SGD, three key differences arise, as discussed below.

595 During early training, the loss may decrease non-monotonically on incorporating momentum, even
 596 at small learning rates. Such oscillations are also observed when quadratic loss functions are
 597 optimized using GD with momentum [11]. These oscillations make it challenging to differentiate
 598 between warmup-induced catapults and fluctuations in loss due to the intrinsic effects of momentum.
 599 Nevertheless, we can still observe loss spikes correlated with an abrupt decrease in sharpness at large
 600 learning rates, as detailed in Appendix E.2.

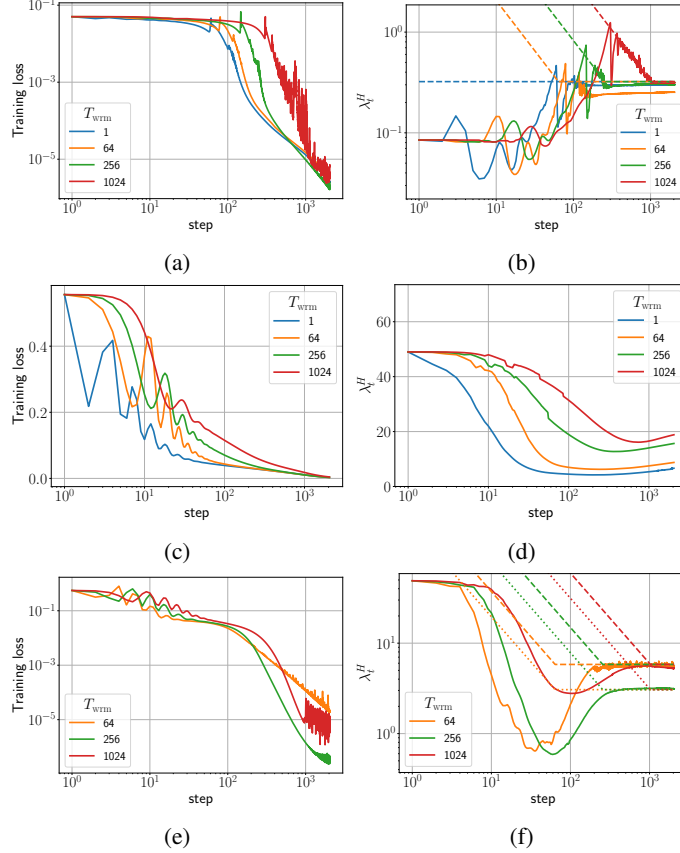


Figure 6: Training loss and sharpness trajectories of FCNs trained on 5k subset of CIFAR-10 using MSE loss and full batch GD with momentum $\beta = 0.9$: (top) μP with learning rate $1/\lambda_0^H$ (middle) SP with learning rate $1/\lambda_0^H$, and (bottom) SP with learning rate $32/\lambda_0^H$. The dotted lines in the sharpness figures correspond to the $(2+2\beta)/\eta_t$ curves, while dashed lines show the $2/\eta_t$ for reference.

601 Additionally, the instability threshold η_c itself evolves differently during training. It changes from
602 $2/\lambda_0^H$ at initialization to $(2+2\beta)/\lambda_t^H$ later in training. Moreover, the late-time instability threshold is
603 significantly influenced by the batch size, exhibiting a much smaller value than SGD for the same
604 batch size. These properties make it more challenging to analyze the training dynamics of SGD with
605 momentum. Nonetheless, the fundamental warmup mechanisms closely mirror the vanilla SGD case.
606 We leave a more detailed analysis of the early training dynamics of SGD-M for future studies.

607 Besides these differences, we note that the warmup mechanisms of SGD with momentum are similar
608 to the vanilla SGD case. We leave a thorough analysis of the early sharpness dynamics of SGD with
609 momentum for future works.

610 E.3 Stochastic Gradient Descent and Cross-entropy Loss

611 The warmup mechanisms for models trained with cross-entropy loss exhibit trends similar to those
612 observed with MSE loss with one crucial difference. Near convergence, sharpness first increases and
613 then abruptly decreases. The decrease in sharpness towards the end of training is observed in previous
614 studies analyzing SGD with fixed learning rate [6]. Additionally, we observe higher fluctuations
615 compared to the MSE loss case. Figure 8 shows trajectories of FCNs under different parameterizations
616 trained on CIFAR-10 with cross-entropy loss using vanilla SGD. Meanwhile, Figure 9 shows the loss
617 and sharpness trajectories of FCNs in SP trained on CIFAR-10 with cross-entropy loss using full
618 batch GD with and without momentum.

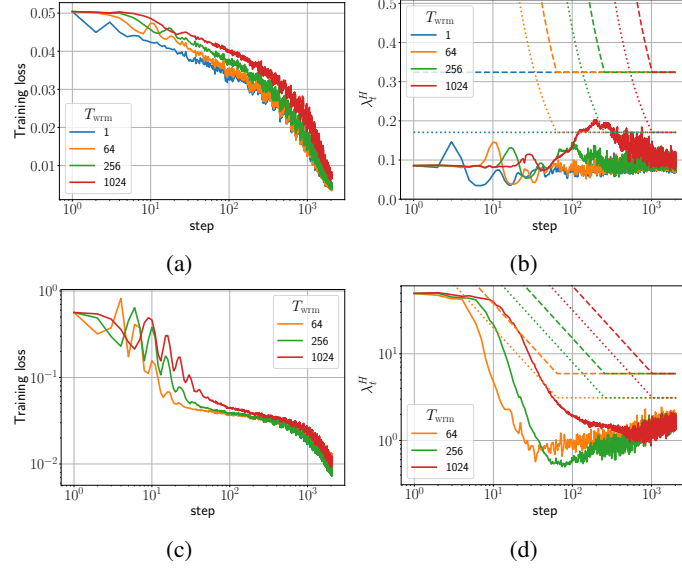


Figure 7: Training loss and sharpness trajectories of FCNs trained on CIFAR-10 with MSE loss using SGD with a batch size $B = 512$ and momentum $\beta = 0.9$: (top) μP with learning rate $1/\lambda_0^H$, and (bottom) SP with learning rate $32/\lambda_0^H$. The dotted lines in the sharpness figures correspond to the $(2+2\beta)/\eta_t$ curves, while dashed lines show the $2/\eta_t$ for reference. Similar mechanisms are observed for cross-entropy loss with a decrease in sharpness at late training times, as detailed in Appendix E.3.

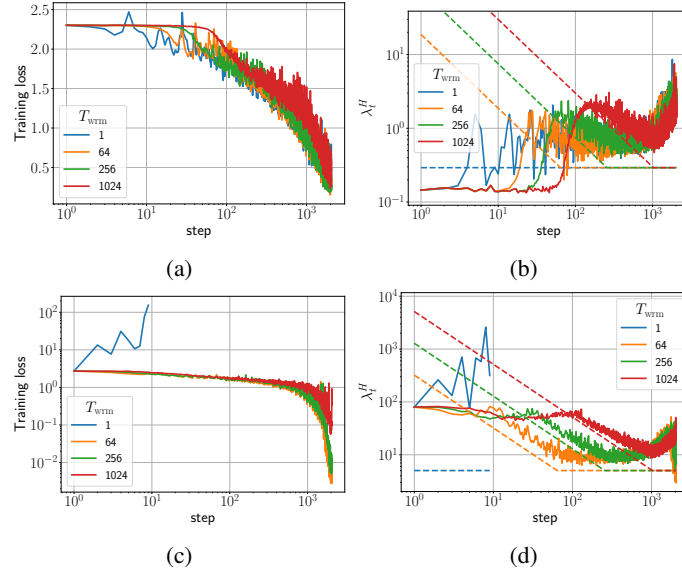


Figure 8: Training loss and sharpness trajectories of FCNs trained on CIFAR-10 with cross-entropy loss using SGD with a batch size $B = 512$. (Top row) μP with learning rate $1/\lambda_0^H$ (Bottom row) SP with learning rate $32/\lambda_0^H$.

619 E.4 Warmup Mechanisms of Adam

620 As discussed in Section 3.2, the instability threshold for Adam is determined by the pre-conditioned
 621 sharpness $\lambda^{P^{-1}H}$ and not by the sharpness itself. Moreover, training dynamics falls under the
 622 sharpness reduction case as the pre-conditioned sharpness starts off large and reduces considerably
 623 during the first few training.

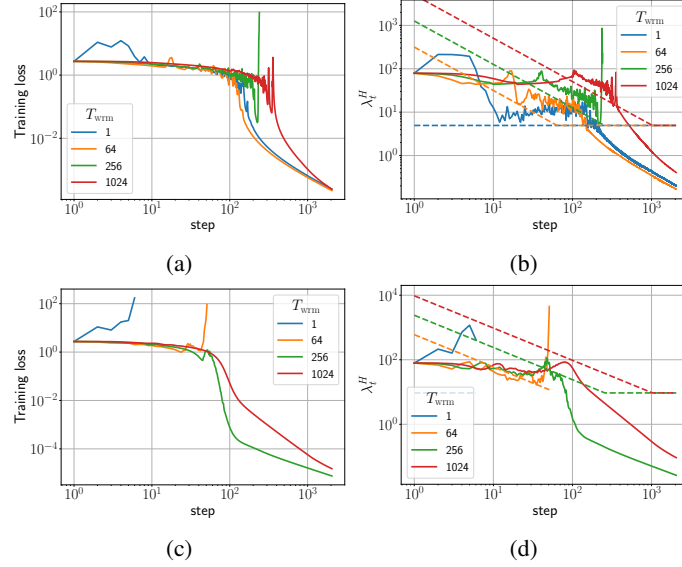


Figure 9: Training loss and sharpness trajectories of FCN-4-512 in SP trained on 5k subset of CIFAR-10 with cross-entropy loss using full batch GD with learning rate $32/\lambda_0^H$ with momentum coefficient (top) $\beta = 0.0$ and (bottom) $\beta = 0.9$.

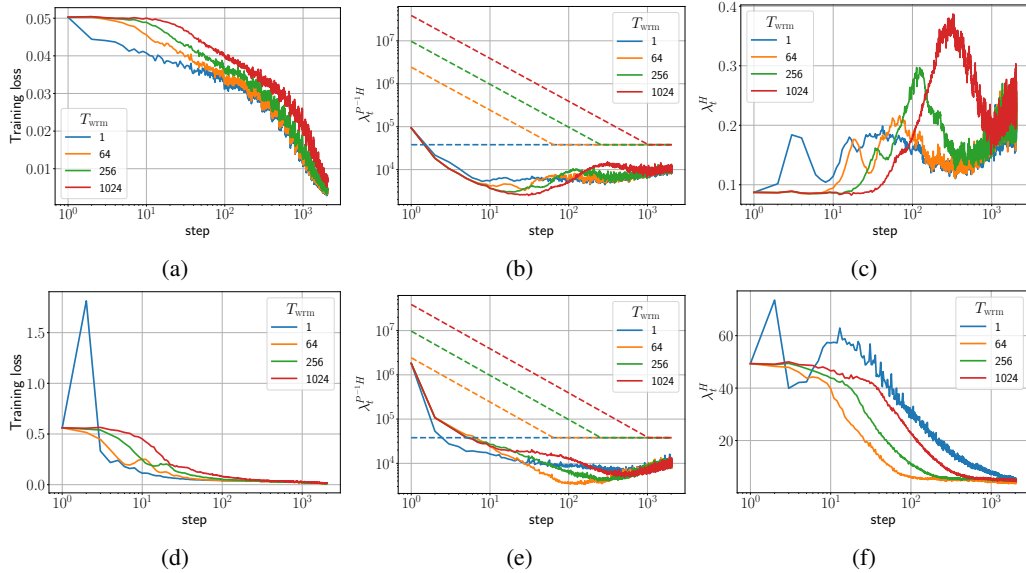


Figure 10: Training loss and sharpness trajectories of FCN-4-512 in (top) μP and (bottom) SP trained on CIFAR-10 with MSE loss using Adam with learning rate $\eta = 0.001$, batch size $B = 512$, $\beta_1 = 0.9$ and $\beta_2 = 0.999$. The dashed lines in the sharpness figures illustrate the instability thresholds $(2+2\beta_1)/\eta_t(1-\beta_1)$.

624 Figure 10 shows the training trajectories of FCNs trained with Adam in the same setting as in
 625 Figure 2 but with a batch size of $B = 512$. Similar to the SGD with momentum case, the late
 626 time sharpness oscillates far below the instability threshold $((2+2\beta_1)/\eta_t(1-\beta_1))$, suggesting that the
 627 instability threshold heavily decreases with a smaller batch size. We note similar findings by Ref. [7].

628 Next, Figure 11 show the warmup mechanism of FCNs trained with cross-entropy loss using Adam
 629 under the full-batch setting. Similar to the SGD case, the pre-conditioned sharpness decreases towards
 630 the end of training.

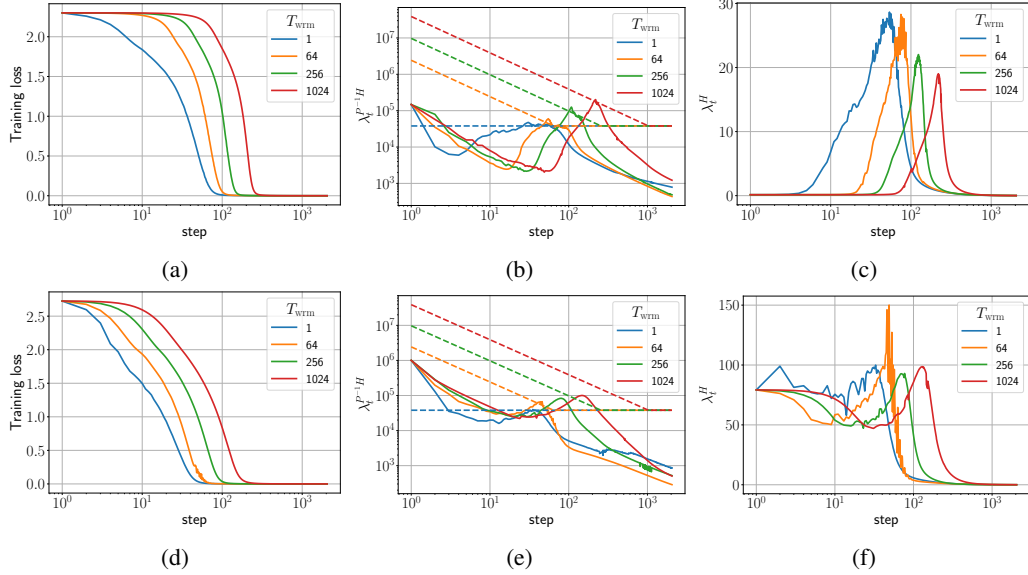


Figure 11: Training loss and sharpness trajectories of FCNs in (top) μP and (bottom) SP trained on CIFAR-10 with cross-entropy loss using full-batch Adam with learning rate $\eta = 0.001$, $\beta_1 = 0.9$ and $\beta_2 = 0.999$. The dashed lines in the sharpness figures illustrate the instability thresholds $(2+2\beta_1)/\eta_t(1-\beta_1)$.

631 E.5 Different Architectures and Datasets

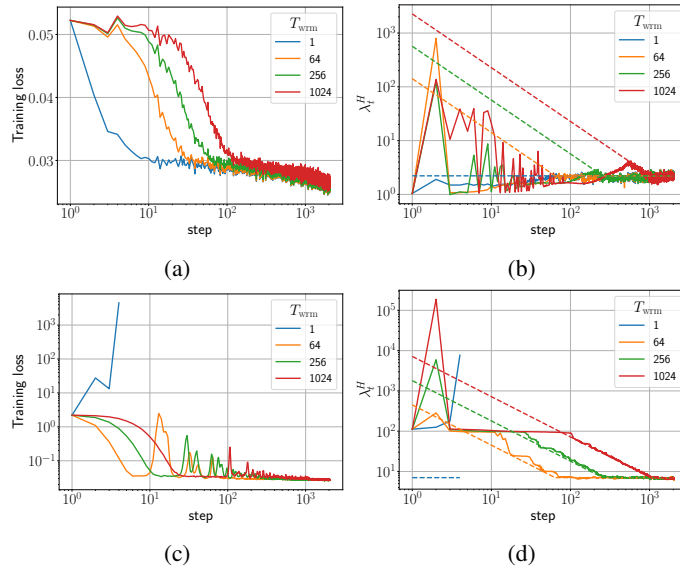


Figure 12: WRN-16-1 trained on CIFAR-10 with MSE loss using vanilla SGD with batch size $B = 512$: (top) μP with $\eta_{\text{trgt}} = 1/\lambda_0^H$ and (bottom) SP with $\eta_{\text{trgt}} = 32/\lambda_0^H$.

632 In the previous sections, we confined our analysis to FCNs to thoroughly explore the effects of
 633 different optimizers and loss functions. This section expands on those results by demonstrating
 634 that the observed warmup mechanisms apply to ResNets and Transformers as well. The Resnet
 635 experiments also employ data augmentation as detailed in Appendix D.1.

636 Figures 12 and 13 show the training trajectories of WideResNets (WRNs) trained on CIFAR-10 with
 637 MSE and cross-entropy loss using SGD. These trajectories generally reflect the warmup mechanisms
 638 discussed in Section 3. However, certain additional features obscure the clarity of these mechanisms.

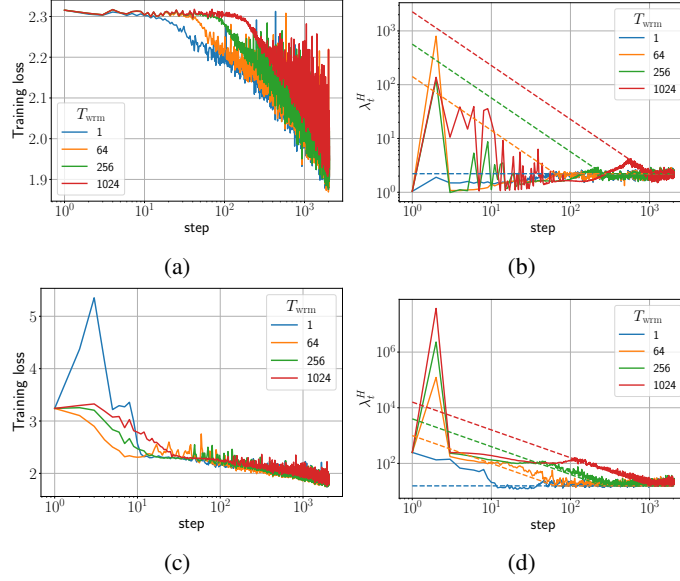


Figure 13: WRN-16-1 trained on CIFAR-10 with cross-entropy loss using vanilla SGD with batch size $B = 512$: (top) μP with $\eta_{\text{tgt}} = 1/\lambda_0^H$ and (bottom) SP with $\eta_{\text{tgt}} = 32/\lambda_0^H$.

639 Notably, we observed a significant sharpness spike on the first training step when using longer
 640 warmup durations, which automatically resolves in the subsequent step. The magnitude of this spike
 641 increases with longer warmup periods. Further analysis revealed that this phenomenon is associated
 642 with an initial increase in the first LayerNorm parameters, which also resolves automatically by the
 643 second step. Beyond this observation, the training trajectories align with the warmup mechanisms
 644 described in the main text.

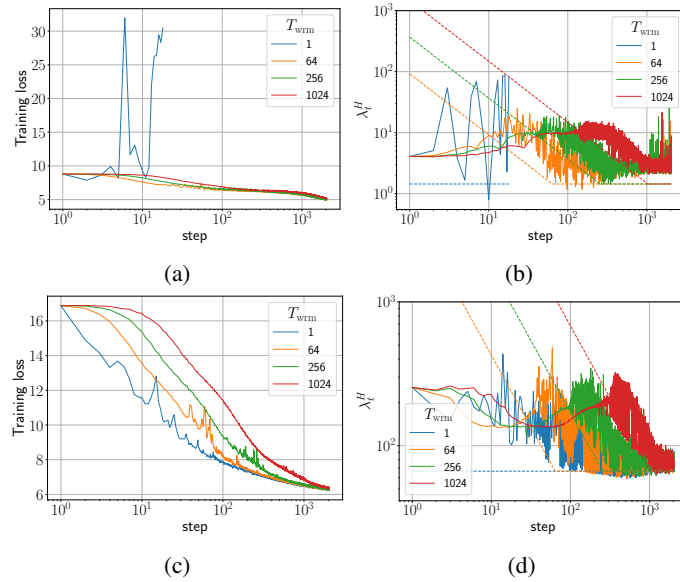


Figure 14: LM-4-128 trained on the WikiText-2 dataset with cross-entropy loss using SGD with a batch size $B = 512$ and a context length $T_{\text{ctx}} = 64$. The top row shows the warmup mechanisms of a Pre-LN Transformer with $\eta_{\text{tgt}} = 5.65/\lambda_0^H$, while the bottom row shows the results for the same Pre-LN Transformer but with the last LayerNorm removed and a learning rate of $\eta_{\text{tgt}} = 8/\lambda_0^H$.

645 Figure 14 illustrates the warmup mechanisms of Pre-LN Transformers trained on the WikiText-2 with
 646 SGD. The Pre-LN Transformer (top row) starts in a flat landscape region ($\lambda_0^H \sim 5$) and experiences

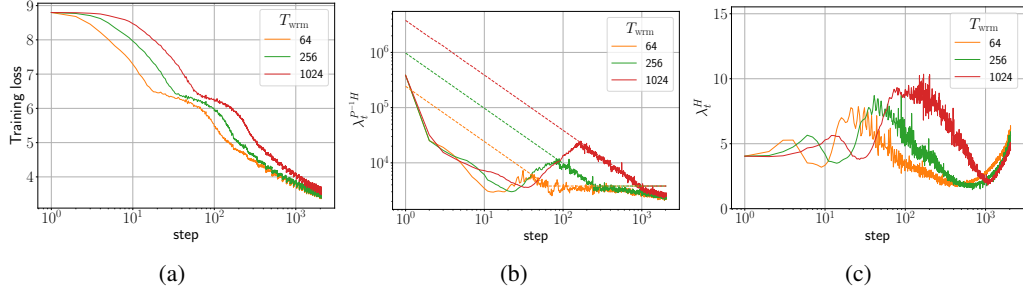


Figure 15: Pre-LN LM-4-128 trained on the WikiText-2 dataset with cross-entropy loss using Adam with a target learning rate $\eta_{\text{tgt}} = 0.003$, a batch size $B = 512$ and a context length $T_{\text{cntx}} = 64$.

647 progressive sharpening right from initialization. In contrast, when the last LayerNorm (just before the
 648 final linear layer) is removed (bottom row), the model starts training in a significantly sharper region,
 649 with the initial sharpness 100 times larger than the standard Pre-LN Transformer. This modified
 650 Pre-LN Transformer experiences a reduction in sharpness during the early stages of training.

651 Figure 15 presents the warmup mechanisms of Pre-LN Transformers trained on WikiText-2 using the
 652 Adam optimizer. Consistent with the results in the main text, the pre-conditioned sharpness exhibits a
 653 reduction early in training, despite the model initializing in a very flat region.

654 These experiments demonstrate that Transformers trained on language modeling tasks exhibit warmup
 655 mechanisms consistent with those discussed in the main text.

656 F Additional Phase Diagrams

657 This section presents further results related to the phase diagrams of warmup shown in Section 4.

658 F.1 Phase Diagrams for different Models and Datasets

659 Figure 16 shows the test accuracy heatmaps of WRN-16-4 trained on CIFAR-100 and Tiny-ImageNet.
 660 These models are trained using cross-entropy loss using SGD with a batch size of $B = 128$.
 661 Additional phase diagrams for Adam are presented in Appendix F.3.

662 Figure 17(a) shows the test loss heatmaps of Pre-LN Transformer trained on the WikiText-2 dataset
 663 using SGD with a batch size $B = 64$. Figure 17(b) shows the Pre-LN Transformer under the same
 664 setup except for the last layer LayerNorm removed. The standard Pre-LN Transformer starts off with
 665 a small sharpness, while the version without the last LN starts off with 100 times higher curvature
 666 and requires warmup to achieve good performance.

667 F.2 The Effect of Momentum and Learning Rate Decay

668 Figure 18 shows that incorporating momentum and cosine decay (for details, see Appendix D.3.3)
 669 minimally affects the warmup phase diagrams. While the conclusions regarding warmup presented in
 670 the main text remain unaffected, we note a few interesting observations.

671 First, the divergent boundary shifts leftward on incorporating momentum, indicating that momentum
 672 permits smaller target learning rates without warmup, and warmup helps SGD-M more. Meanwhile,
 673 cosine decay has a minimal effect on the divergent boundary.

674 Additionally, we observe a performance enhancement by incorporating momentum, especially at
 675 small learning rates. In contrast, a decaying learning rate beyond warmup degrades performance
 676 at small learning rates while improving at higher ones. Finally, incorporating both momentum and
 677 cosine decay leads to further enhancement, indicating a synergistic interaction between the two.

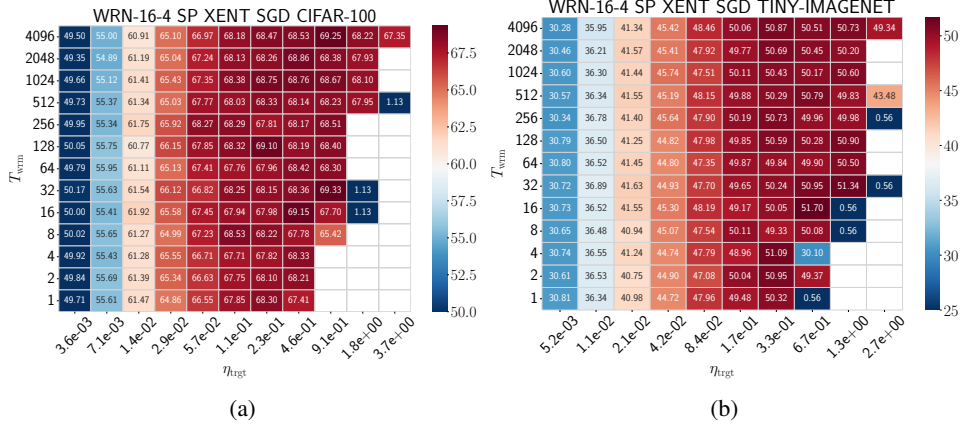


Figure 16: Test accuracy heatmaps of WideResNets (WRNs) in SP trained on (a) CIFAR-100 and (b) Tiny ImageNet with cross-entropy loss using SGD with batch size $B = 128$.

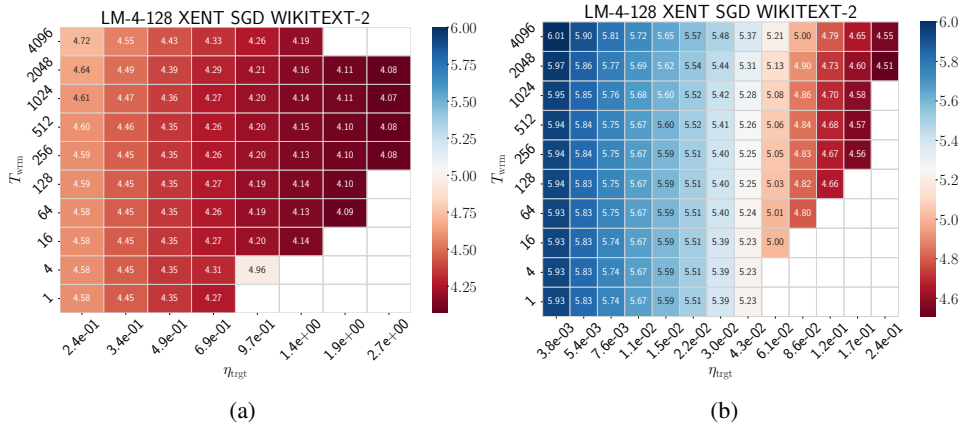


Figure 17: Test loss heatmaps of LM-4-128 in SP trained on WikiText-2 with cross-entropy loss using SGD with a batch size $B = 64$: (a) Pre-LN Transformer and (b) Pre-LN Transformer without the last LayerNorm.

678 F.3 Phase Diagrams of Adam and GI-Adam

679 Figures 20 to 22 compare the warmup phase diagrams of Adam and GI-Adam of WRNs trained on
 680 CIFAR-100, Tiny-ImageNet and of Transformers trained on WikiText-2 dataset. Similar to the results
 681 shown in the main text, GI-Adam enhances performance over standard Adam by pushing the failure
 682 boundary.

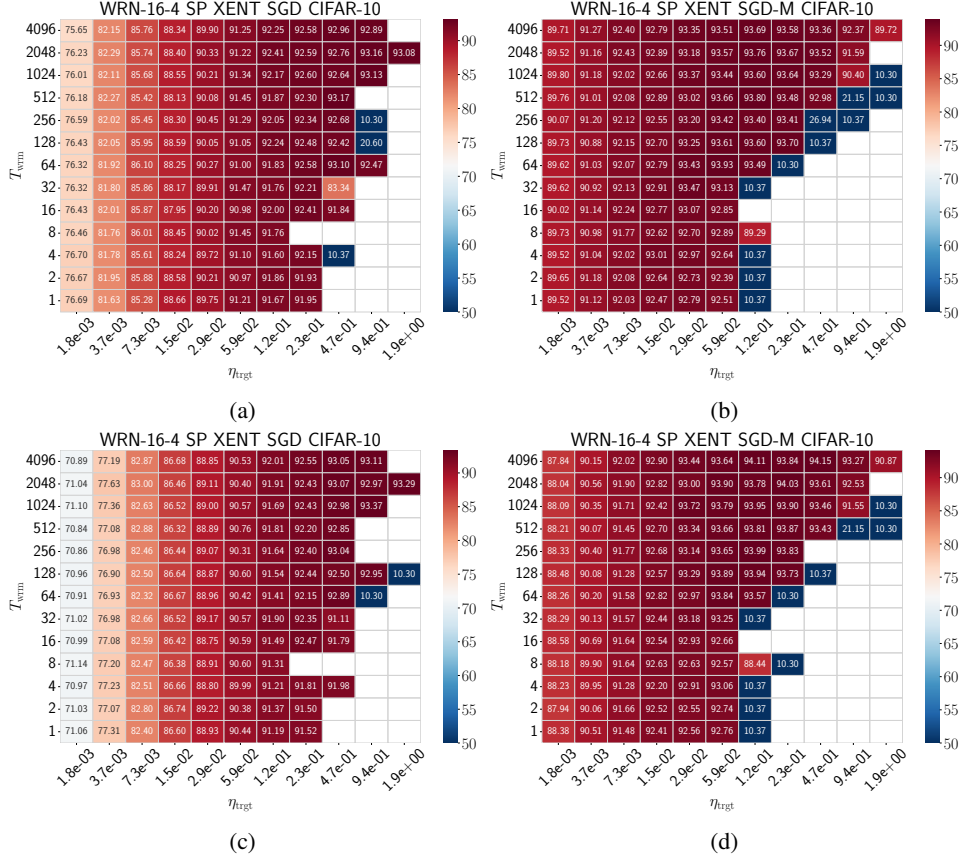


Figure 18: Test accuracy heatmaps of WideResNets (WRNs) in SP trained on CIFAR-10 with cross-entropy loss using SGD with batch size $B = 128$: (top row) no cosine decay (a) no momentum, (b) momentum with $\beta = 0.9$, and (bottom row) with cosine decay (c) no momentum, and (d) momentum with $\beta = 0.9$. The setting of (a) is the same as in Figure 3(c) but with a different mini-batch sequence.

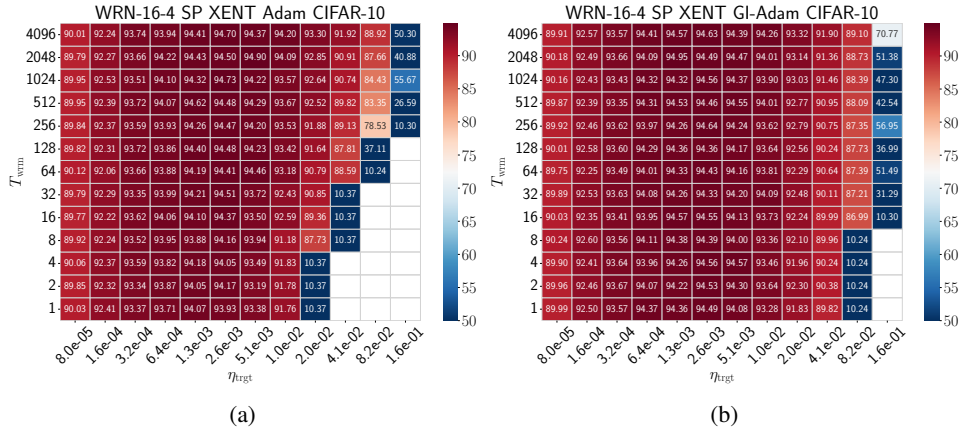


Figure 19: Test accuracy heatmaps of WRN-16-4 trained on CIFAR-100 with cross-entropy loss using (left) standard Adam, and (right) GI-Adam with batch size $B = 128$.

683 G Non-divergence of Adam

684 Figure 23 shows that, despite experiencing catastrophic instabilities during early training, Adam
 685 does not diverge well beyond the training failure boundary. While Adam can recover from these

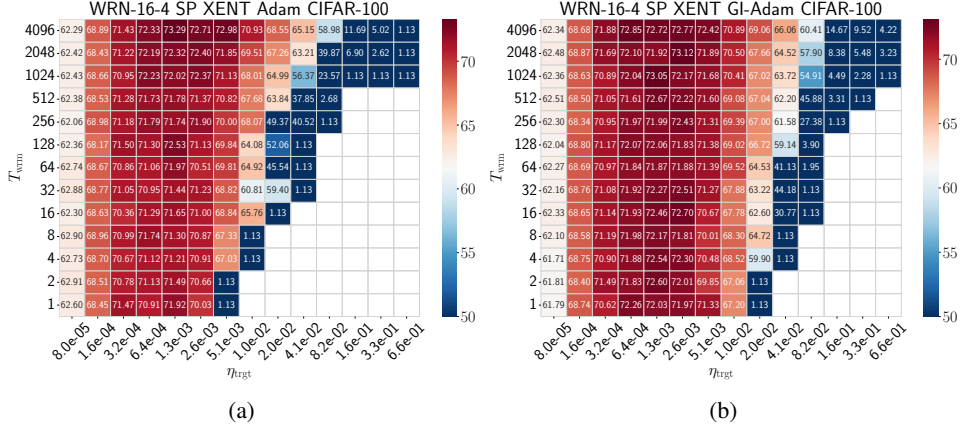


Figure 20: Test accuracy heatmaps of WRN-16-4 trained on CIFAR-100 with cross-entropy loss using (left) standard Adam, and (right) GI-Adam with batch size $B = 128$.

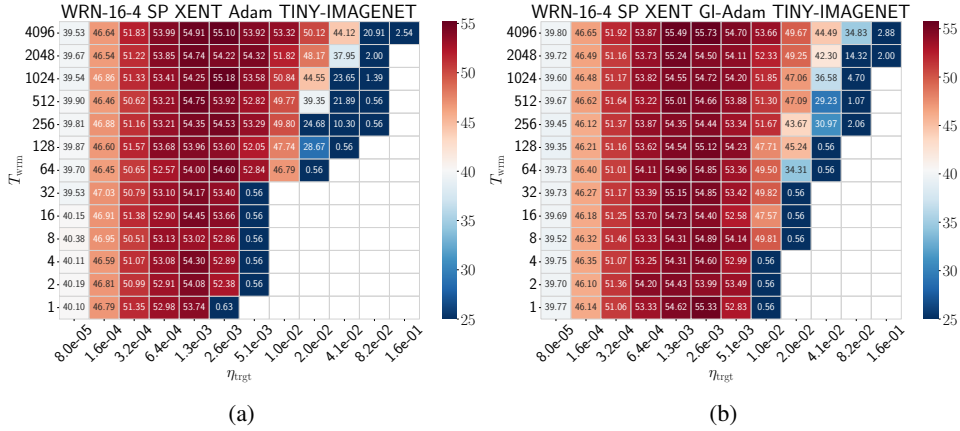


Figure 21: Test accuracy heatmaps of WRN-16-4 trained on Tiny-ImageNet with cross-entropy loss using (left) standard Adam, and (right) GI-Adam with batch size $B = 128$.

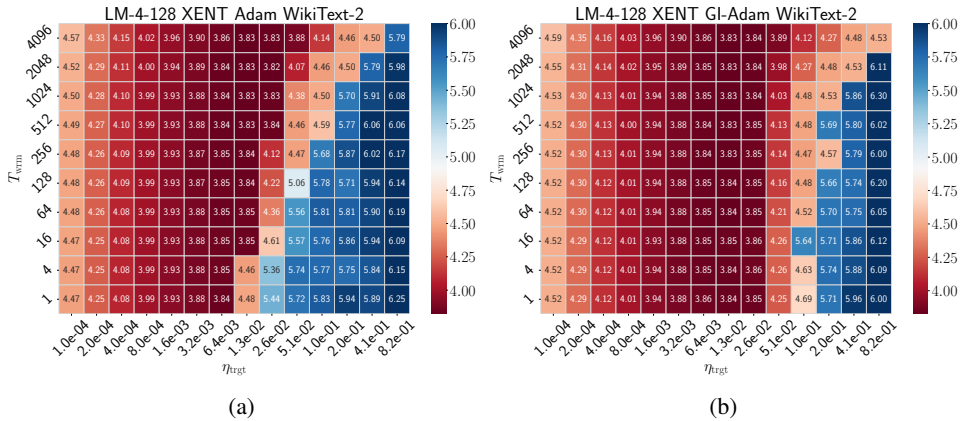


Figure 22: Test loss heatmaps of LM-4-128 in SP trained on WikiText-2 with cross-entropy loss using (a) standard Adam, and (right) GI-Adam with batch size $B = 64$.

686 instabilities, the model’s performance is severely impacted, resulting in training failures rather than
 687 convergence to a reasonable minimum.

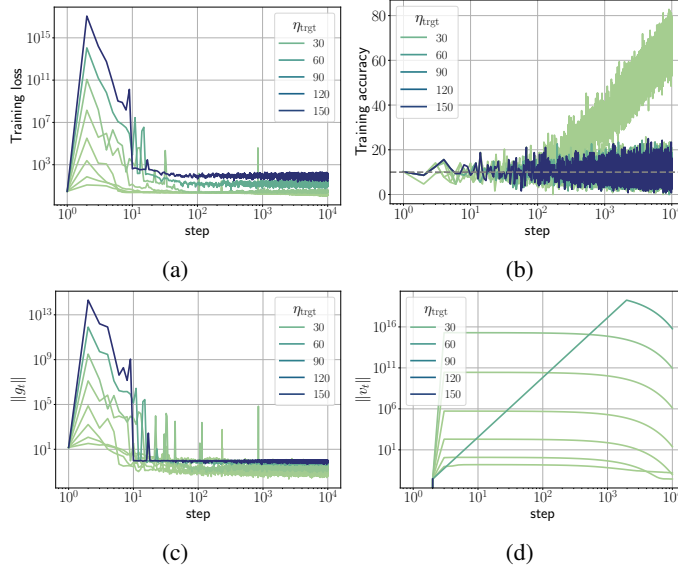


Figure 23: Training trajectories of WRNs trained on CIFAR-10 using Adam with cross-entropy loss and varying learning rates. The setup is identical to the $T_{\text{wrm}} = 1$ row of Figure 4, but without employing cosine learning rate decay. The first training failure is observed at a learning rate of $\eta_{\text{trgt}} = 0.02048$. To investigate the behavior beyond the training failure boundary, learning rates are sampled from $\eta_{\text{trgt}} = 0.01024$ (just below the failure boundary) up to $\eta_{\text{trgt}} \approx 150$.

688 These large loss catapults cause the gradients \mathbf{g} to spike during early training, leading to a substantial
 689 increase in its second moment \mathbf{v} . While the gradients return to a lower value after a few training
 690 steps, the second moment remains large in magnitude for a prolonged period. These large values of
 691 \mathbf{v} result in a small effective learning rate, which hinders training to escape these high-loss regions.
 692 Consequently, the models remain stuck in a suboptimal state rather than converging. We refer to this
 693 as a training failure.

694 Upon closer examination of the individual layers during training failures, we found that certain layers
 695 or residual blocks output zero. This results in vanishing gradients except for the last layer bias and
 696 training halts. We defer the detailed analysis of Adam’s failures to future work.

697 H Additional Results on GI-Adam

698 This section presents additional results for GI-Adam. We provide further insights into the mechanisms
 699 and interpretations of GI-Adam.

700 H.1 Warmup Mechanisms of GI-Adam

701 Figure 24 shows the training trajectories of FCNs with different parameterizations trained with
 702 GI-Adam. Notably, the pre-conditioned sharpness starts at significantly lower values than standard
 703 Adam. Specifically, for the μP model, the initial pre-conditioned sharpness $\lambda^{P-1}H$ is around 2000
 704 instead of the value 10^5 observed for Adam (c.f. Figure 2). Remarkably, this almost eliminates initial
 705 sharpness reduction. Similarly, the pre-conditioned sharpness for the SP model starts around 10^4
 706 instead of 10^6 . Notably, in the SP scenario, there is no initial spike in the $T_{\text{wrm}} = 1$ (c.f. Figure 2),
 707 demonstrating that this simple modification effectively reduces instabilities during the early training.

708 H.2 GI-Adam as an Automated Warmup

709 In this section, we show that a bias correction is not required when the second moment is initialized
 710 with the gradients at initialization in GI-Adam. Therefore, employing a bias correction as in the
 711 original Adam algorithm in this case serves as an automated warmup given by $\eta_t = \eta_{\text{trgt}} \sqrt{1 - \beta_2^t}$.

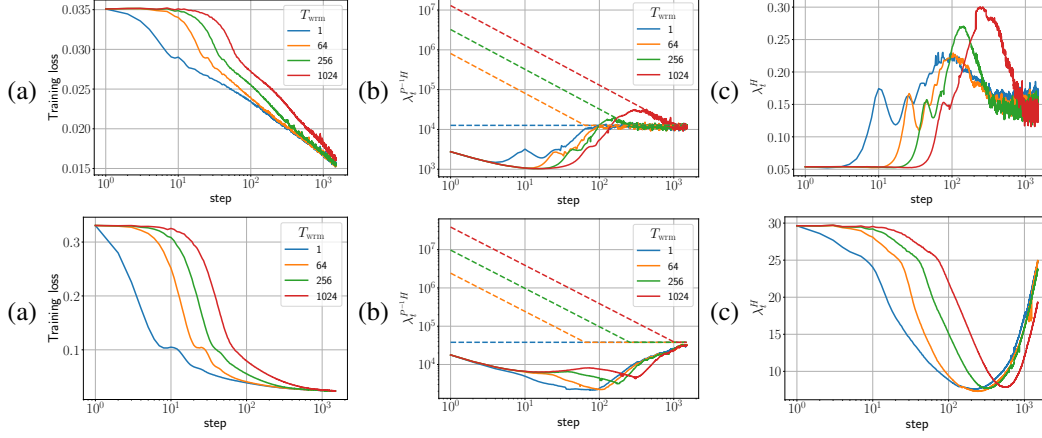


Figure 24: Training loss and sharpness trajectories of FCNs in (top) μP and (bottom) SP . The experimental setup is identical to Figure 2 but with GI-Adam instead of standard Adam.

712 The moving average of the second moment is given by:

$$\mathbf{v}_t = (1 - \beta_2) \sum_{i=0}^{t-1} \beta_2^i \mathbf{g}_{t-i}^2 + \beta_2^t \mathbf{v}_0, \quad (9)$$

713 where $\mathbf{v}_0 = \mathbf{g}_0^2$. Following standard assumptions, we assume that the second moment of the gradient
 714 is constant during early training $\mathbb{E}[\mathbf{g}_t^2] = \sigma^2$. Taking the expectation of the above equation over the
 715 gradient distribution yields

$$\mathbb{E}[\mathbf{v}_t] = (1 - \beta_2) \sum_{i=0}^{t-1} \beta_2^i \mathbb{E}[\mathbf{g}_{t-i}^2] + \beta_2^t \mathbb{E}[\mathbf{v}_0]. \quad (10)$$

716 Simplifying the above equation, we have

$$\mathbb{E}[\mathbf{v}_t] = (1 - \beta_2) \sigma^2 \frac{1 - \beta_2^t}{1 - \beta_2} + \beta_2^t \sigma^2 = \sigma^2. \quad (11)$$

717 This result demonstrates that when the second moment is initialized with the gradients at initialization,
 718 it does not require bias correction, as the expected value of the second moment is equal to the constant
 719 σ^2 . If we apply the usual bias correction on top of initializing the second moment with the gradients,
 720 we effectively downscale the second moment by a factor $\sqrt{1 - \beta_2^t}$. Assuming small enough ϵ , this
 721 can be viewed as a multiplicative factor to the learning rate. As a result, GI-Adam is equivalent to
 722 having a natural warmup given by $\eta_t = \eta_{\text{tgt}} \sqrt{1 - \beta_2^t}$.

723 H.3 The Primary benefit of GI-Adam results from the magnitude of the second moment at 724 initialization

725 To further assess if the primary cause of instability during early training is the large $\lambda^{P-1}H$, we
 726 randomly initialize \mathbf{v}_0 but with the same norm as the gradients at initialization. We refer to this as
 727 Randomly Initialized Adam (RI-Adam). Like GI-Adam, this also results in improved performance as
 728 shown in Figure 25.

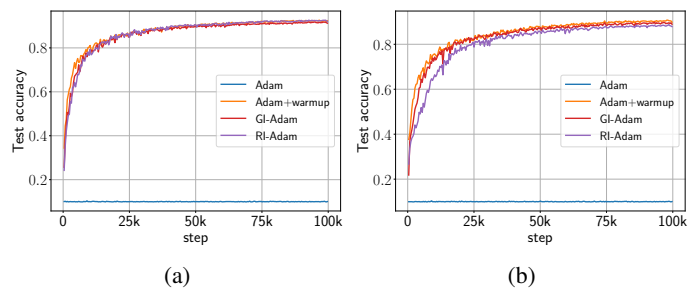


Figure 25: Comparison of test accuracy trajectories of WRNs trained with different Adam variants for two target different learning rates: (a) $\eta_{\text{trgt}} = 0.020480$, and (b) $\eta_{\text{trgt}} = 0.040960$. For Adam+warmup, the warmup duration is set to $T_{\text{wrm}} = 1024$.