Benchmarking Optimizers for Large Language Model Pretraining

Anonymous Author(s)

Affiliation Address email

Abstract

The recent development of Large Language Models (LLMs) has been accompanied by an effervescence of novel ideas and methods to better optimize the loss of deep learning models. Claims from those methods are myriad: from faster convergence to removing reliance on certain hyperparameters. However, the diverse experimental protocols used to validate these claims make direct comparisons between methods challenging. This study presents a comprehensive evaluation of recent optimization techniques across standardized LLM pre-training scenarios, systematically varying model size, batch size, and training duration. Through careful tuning of each method, we provide guidance to practitioners on which optimizer is best suited for each scenario. For researchers, our work highlights promising directions for future optimization research. Finally, by releasing our code and making all experiments fully reproducible, we hope our efforts can help the development and rigorous benchmarking of future methods.

4 1 Introduction

2

3

6

8

9

10

11

12

13

- Over the past five years, Large Language Models (LLMs) [15, 59, 22, 48] have shown growth in performance and size, demonstrating proficiency in various downstream tasks [80, 7, 85]. The success of LLM pretraining hinges on three key pillars: high-quality data [65, 44], architectural innovations [31] [15], and scalable optimization techniques.
- Among these, the choice of optimizer has remained notably consistent in recent years, with Adam(W) [38, 50] dominating deep learning for nearly a decade. However, recent advances [33, 47, 84, 62, 66, 17] challenge this status quo, offering alternatives that surpass AdamW in speed, communication efficiency [1] or final downstream performance on various benchmarks [12, 37], particularly for autoregressive language modeling [70]. Despite these innovations, current benchmarks and ablation studies [96, 34] remain narrow in scope, often examining only isolated aspects of optimizer design. This lack of systematic comparison makes it difficult to obtain trustworthy insights for practitioners, or identify the next promising research directions.
- In this work, our goal to revisit the problem of benchmarking optimizers for LLM pretraining.
 We do so through standardized experiments which vary important parameters such as batch size,
 model size, and the number of training iterations. This allows us to formulate an up-to-date list of
 best-performing methods for the community of researchers and practitioners. We demonstrate the
 efficiency of each considered method through careful tuning, and present insightful ablations along
 the way. Furthermore, we provide a set of best practices for LLM pretraining that are applicable
 regardless of the optimizer chosen.
- We summarize our contributions as follows:

(Contribution 1) We conduct the first large-scale, controlled benchmark of 11 different optimization
 methods across diverse LLM training scenarios. A fair comparison is ensured by precise accounting
 for compute costs, and extensive hyperparameter tuning. We identify optimal optimizer choices in
 several relevant training regimes, for both dense and MoE architectures.

(Contribution 2) We perform comprehensive ablations of critical training hyperparameters—including warmup duration, initialization schemes, gradient clipping, final learning rates, and learning rate scheduler choices—providing actionable insights for optimizing LLM training in practice.

(Contribution 3) We open-source our full benchmarking toolkit, including training scripts, evaluation pipelines, and hyperparameter configurations, to enable reproducible research and facilitate future optimizer development.

For practitioners, our work provides an evidence-45 based answer to the burning question: "Is Adam still 46 the most effective optimizer in the age of LLMs, or can we achieve better performance at scale with 48 novel optimizers?". For researchers, our work de-49 livers a unified benchmarking framework for LLM 50 pretraining, along with extensive ablation studies 51 which systematically evaluate both popular and over-52 looked optimizer designs—revealing previously un-53 explored tradeoffs between efficiency, stability, and 54 final model performance. Overall, our findings not only challenge long-held assumptions about opti-56 mizer selection but also establish a foundation for 57 future advances in large-scale model training. By 58 bridging the gap between theoretical innovation and 59 practical deployment, this work aims to accelerate 60 progress in both research and industry applications 61 of LLM training.

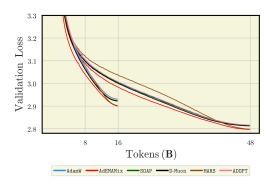


Figure 1: A comparison of leading optimizers, for training a 720M parameter LLM.

63 2 Background & Related Work

Optimizers. While computer vision models often show comparable performance between SGD [72] and AdamW [94], the landscape differs dramatically in LLM training. Recent work [95] demon-65 strates that adaptive methods like AdamW provide substantially better optimization characteristics for transformer-based language models. The question of why AdamW works so well has been a long-standing topic of research [2, 60, 93, 43, 41]. Modern methods often inherit AdamW's core ideas in their structure, such as ADOPT [83] and AdEMAMix [62]. ADOPT has been motivated by solving 69 long-standing convergence issues in AdamW. By normalizing the second-order moment prior to the 70 momentum update, they eliminate the non-convergence issues of AdamW on smooth non-convex 71 functions. Meanwhile AdEMAMix extends AdamW with an additional slower momentum buffer, i.e. a 72 slower exponential moving average (EMA), which allows the use of much larger momentum values, 73 accelerating convergence. 74

One interpretation of AdamW's effectiveness lies in its sign-based update [42]: without the exponential moving average (EMA), AdamW resembles signSGD [6]. Recent work [96, 36] has shown that Signum (signSGD with momentum), can perform comparably to AdamW. Earlier, the community also discussed Lion [9], a method with a similar sign-based structure. Signum and Lion offer memory benefits due to the use of only a single instead of Adam's two buffers for optimizer states.

Another family of methods stems from AdamW's approximate second-order structure, where the diagonal of the Fisher information matrix or other preconditioning approaches [52, 24] are used as the second moment estimate. This idea has given rise to Sophia [46], SOAP [84], and, to some extent, Muon [33].

The parameter-free concept [61] has led to the development of Schedule-Free AdamW (SF-AdamW) [17] and Prodigy [54]. These optimizers do not require a decreasing learning rate schedule, making them relevant for continual training. Last but not least, MARS [88], builds upon this line of research and incorporates a variance reduction mechanism in its update rule.

Benchmarks. To a large extent, the benchmarking setup determines the final conclusions. Some benchmarks are designed for short speedruns in terms of training or validation loss [32], while 89 others focus on a downstream target metric after training [96, 12, 76]. Methods that perform well 90 in short speedruns might not be optimal for longer training horizons as in real LLM training runs 91 (see Figure 3). "But what constitutes a sufficiently long horizon?" "What should be the compute 92 budget for LLM training?" These are questions explored by scaling laws [35]. Early benchmarks 93 for optimizers and other ablation studies often rely on Chinchilla scaling laws [26] with a ratio of roughly 20 tokens per parameter (TPR) needed for pretraining. However, recent research [69, 74] 95 argues that this is far from sufficient for production-ready models. 96

Another important issue is the choice of loss function. Recent setups have been using an auxiliary z-loss [86] [11] in addition to cross-entropy, which requires further investigation. We believe this choice is influenced by the use of the OLMo [58] codebase, which we also address in our work.

Additionally, we found that previous setups for comparing optimizers do not align with recent best practices regarding weight decay, learning rate decay, and overall hyperparameter tuning. All of these questions are revisited in our work.

3 Experimental Setup

103

125

126

127

128

129

130

131

133

134

135

136

137

138

139

Notations. We use the following notations. Let γ be the learning rate, λ the weight decay coefficient, and T the total number of iterations. Momentum-related parameters are represented by the symbol β .

Models & Data. For most experiments, we use a Llama-like transformer [48] architecture, including SwiGLU activations [77], RMSNorm [91], and RoPE embeddings [81]. We experiment with four sizes of models: 124M, 210M, 583M, 720M. We train on a 100B tokens subset of FineWeb [64]. It consists of a cleaned and deduplicated corpus for LLM pretraining, which we tokenize using the GPT-2 tokenizer prior to splitting into train and validation sequences. MoE setup described in

Iterations & Batch size. Throughout our experiments, we use a sequence length of 512 tokens. For 111 clarity, we often report the batch size in tokens by writing *Batch size* \times *sequence length*. For the 124M 112 model, we use batch sizes of $32 \times 512 = 16$ **k**, $256 \times 512 = 131$ **k**, and $512 \times 512 = 262$ **k** tokens; 113 for the 210M model, we use a batch size of $256 \times 512 = 131$ k; and for 583M model, we leverage 114 the batch sizes of $1024 \times 512 = 524$ k and $3936 \times 512 = 2$ M tokens. Depending on the model size, 115 we vary the number of iterations — also measured in tokens for compatibility with scaling laws and 116 to accommodate different batch size settings. We train 124M and 210M models for equal durations 117 of $\{1, 2.1, 4.2, 6.3, 8.4, 16.8\}$ **B** tokens. This corresponds to $T \in \{64, 128, 256, 384, 512, 1024\}$ **k** 118 iterations for a batch size of 32, and $T \in \{8, 16, 32, 48, 64, 128\}$ **k** iterations for a batch size of 256. For 583M models, we train on $\{13,32\}$ B tokens, corresponding to $T \in \{6.5,16\}$ k iterations, resp. $T \in \{25, 61.5\}$ k iterations, for a batch size of 3936, resp. 1024. In the setup with 720M model, 121 we have $T \in \{8, 16, 48\}$ k iterations for a batch size of 1M tokens. Thus, for all model scales, 122 we include both Chinchilla optimal lengths of training and beyond. More details are available in 123 Appendix C. 124

Loss. We train using the classical cross-entropy next token prediction loss. Some prior works introducing optimizers use a z-loss in addition to cross-entropy [30, [11, 86, 84, 96]]. We found that this has little impact and, therefore, do not use z-loss. An ablation showing results with and without z-loss is in the appendix.

Hyperparameter Tuning. Training LLMs is a computationally intensive task. As a guidance, practitioners often rely on insights gathered at lower scales, scaling laws, and other rules [87] [18]. It is also commonplace to run experiments for only a shorter duration of training, as a way to test certain hyperparameters prior to extending the training horizon to more iterations. Because a full grid search over every hyperparameter, for each setting and optimizer, would be too costly, we resort to a similar approach. More precisely, for each model size, batch size, and optimizer, we tune optimization hyperparameters extensively for a number of training tokens which is near-Chinchilla optimal. We then keep those hyperparameters when we increase the number of iterations. While we found that the sensitivity to several hyperparameters can change as we increase the training horizon, we found this approach simple and yet effective. The hyperparameters being considered depend on the optimizer. We proceeded from small to large model scale, and used insights gathered at smaller scales

https://huggingface.co/datasets/HuggingFaceFW/fineweb

to guide the hyperparameter search at larger scales. Our hyperparameter sweeps are summarized in Appendix D. We present the clarifications regarding the connection between the number of iterations 141 and tokens for different batch size settings as well as the Chinchilla optimal length of training for 142 our models in Tables 3 and 5. As learning rate schedulers, we compare cosine [49], linear and 143 warmup-stable-decay (WSD) [27, 90, 28]. Unless specified, we use a cosine scheduler. Results with 144 WSD and linear schedulers are discussed in Section 4. Recent works also emphasize the importance 145 146 of sufficiently decaying the learning rate [4, 75] [28]. As such, we take care to decay to $0.01 \times \gamma$ instead of the often used $0.1 \times \gamma$. To give an idea of how much effort was put into tuning each 147 method, across all model sizes, batches and iterations, we trained a total of 2400 models, and have 148 spent roughly 30000 GPU hours. 149

Optimizers. Here is a list of the optimizers we considered in our work. For each algorithm, we 150 write in parentheses the optimizer-specific hyperparameters we tuned: AdamW(β_1, β_2), SOAP(β_1, β_2) 151 and preconditioning frequency, Lion (β_1, β_2) , MARS (η, β_1, β_2) and Newton-Schulz hyperpa-152 rameters, ADOPT(β_1, β_2), Signum(β), Prodigy(β_1, β_2), SF-AdamW(β_1, β_2), Muon($\gamma^{M}, \beta, \beta_1, \beta_2$), Sophia (ρ, β_1, β_2) , AdEMAMix $(\beta_1, \beta_2, \beta_3, \alpha)$. When an optimizer has several momentum variants 154 e.g. Nesterov [57] or Polyak [67], we try both. In addition, we tune the learning rate γ extensively 155 for all methods. We also try different gradient clipping, warmup steps, and weight-decay values. A 156 summary of the hyperparameters tested and selected for each model size is in Appendix D. All the 157 optimizers are described in depth in Appendix A 158

Results

159

163

164

165

167

168

169

170

171

172

173

174

175

176

181

183

184

186

187

188

189

190

191

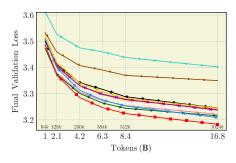
We structure our story starting with smaller models and batch sizes, and gradually scaling up to larger 160 configurations. In some instances, we complement the core benchmarking results with additional 161 ablations and possible best-practices. 162

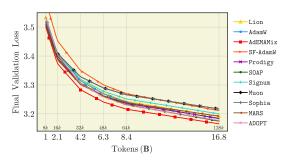
Benchmarking at Small Scale: Training Models of 124M Parameters

Using "small" batches. We first report results when using batches of 32×512 tokens in Figure 3. We tune hyperparameters by training for 2.1B tokens (128k iterations) and then keep those hyperparameters for all other training durations. The best hyperparameters are reported in Appendix D.I. 166 We observe how, for the smallest number of iterations we considered (1B tokens $\equiv 64$ k), SOAP, ADOPT and AdEMAMix all outperform AdamW, with SOAP being the best. As we increase the number of iterations, AdEMAMix takes the lead while AdamW closes the gap with both ADOPT and SOAP. A sign-based methods such as Lion and Signum are expected to perform poorly when the batch size is small. Intuitively, this is due to the sign (\cdot) operator being sensitive to gradient noise. As described in its original paper, MARS also performs poorly when the batch size is small. We found Prodigy, Muon and SF-AdamW to underperform in this setting compared to AdamW. On this scale, Prodigy suffers from the lack of bias correction of the learning rate, as well as being sensitive to (β_1, β_2) (see Figure 17

Using "large" batches. We now report results when using batches of 256×512 tokens — $8 \times$ larger than for our small batch setting. Results in Figure 2 show how Signum, Mars, Lion, Prodigy greatly benefit from the increased batch size. Remarkably, we observe that the Prodigy method scales similarly to AdamW. We emphasize the possible community interest in this algorithm as its EMA Prodigy adaptively emulates the learning rate behaviour. For a small number of iterations 180 (e.g. $T \in \{8\mathbf{k}, 16\mathbf{k}\}$ corresponding to 1B and 2B tokens), all methods outperform AdamW except for SF-AdamW and Sophia. As we increase the number of iterations ADOPT, SOAP, and AdEMAMix take 182 the lead. In particular, AdEMAMix has a consistent lead over other methods. While we anticipated—in accordance with Vyas et al. [84]—that SOAP would greatly benefit from the larger batch size, its behavior remains relatively consistent compared to our previous small batch setting.

Stability across training horizons. As mentioned in Section 3, we tune hyperparameters training on 2.1B tokens and keep those hyperparameters when extending the training horizon. In Figure 3 we study whether it is possible to find better parameters for AdamW, SOAP, and AdEMAMix. When training on 16.8B tokens, we see it is beneficial to increase the β_3 from 0.999 to 0.9999. Without this improvement, SOAP ends up matching the performances of AdEMAMix when extending the training horizon further to 33.6B tokens ($\equiv 256k$ iterations). In our experiments, $\beta_3 = 0.999$ is only better than $\beta_3 = 0.9999$ when the number of training iterations is less than 32k. This matches observation





(a) Batch size 32×512 tokens.

(b) Batch size 256×512 tokens.

Figure 2: Ranking of optimizers for 124M models with small and large batch sizes. In both (a) and (b), we show the *final* validation loss for different training durations, corresponding to different numbers of tokens. Above each token number, we write the number of training iterations corresponding. In (a), we use a "small" batch size of 32×512 tokens. In (b), we use a larger batch size of 256×512 tokens.

from [62], which recommends reducing β_3 when training for fewer iterations. We also test whether the learning rate γ changes as we increase the number of tokens/iterations. In Figure [5], we run a sweep over γ when training for $16.8\mathbf{B}$ tokens. While for most methods, the best γ obtained in the previous sweep remains optimal, this is not the case for SOAP and SF-AdamW, which can benefit from a larger $\gamma = 0.002$.

WSD vs. cosine & linear γ -schedulers. Learning rate schedulers received a lot of attention recently [79, 28]. We conducted a series of experiments comparing WSD [27, 90] and linear with cosine [49] learning rate schedulers. Surprisingly, the performance gap between these two schedulers observed in Figure [23] is often significant [2] for benchmarking optimizers. Consequently, we decided to adopt the cosine scheduler for all further experiments.

Decaying γ sufficiently. In Figure 8 we show the impact of decaying more or less the learning rate $\gamma^{(t)}$. From $\gamma=10^{-3}$ we train models using cosine decay down to $\gamma_{\rm end} \in \{10^{-4}, 10^{-5}, \dots, 10^{-9}\}$. We found that decaying the learning rate sufficiently matters. In particular, the often use rule consisting in decaying to $0.1 \times \gamma$ is suboptimal. This agrees with the recent works [28], [75], [4]. Building on this findings, we consistently use cosine decay down to $0.01 \times \gamma$.

Takeaway 1. After the experiment in the small-batch setting, we conclude that: (i) AdEMAMix scales in the best manner with the number of iterations, SOAP underperforms AdamW when the length of training increases. ADOPT and Prodigy show almost equal performance across all training durations. Sign-based methods, predictably, underperform when the batch size is small, but what is interesting, is that Sophia diverges at all, even if trains with sufficiently small learning rate.

Increasing the batch size further. We also run an experiment with batches of $512 \times 512 = 262 \mathbf{k}$ tokens, training for $64 \mathbf{k}$ iterations. Results in Figure 3 show mostly consistent results. Noticeably MARS becomes the second best performing method behind AdEMAMix, followed closely by Prodigy, Lion, and SOAP. Interestingly, Signum performs comparably to AdamW.

Takeaway 2. Taking into consideration large batch size setting, we found that many methods, once properly tuned, can show a remarkable performance compared to AdamW and also outperform it.

Weight decay ablation. As recent frameworks for LLM pretraining or ablation studies omit weight decay as a default non-zero hyperparameter, some setups even mislead by not incorporating weight decay in their experiments. In this work, we demonstrate the importance of weight decay and its impact across different optimizers. Surprisingly, increasing weight decay while keeping the learning

²We emphasize that the difference between the two schedulers is generally less than 5% of the total compute spent. However, this still represents a significant gap in our benchmarking setup, e.g., SF-AdamW may outperform AdamW in some settings (see Figure 23).

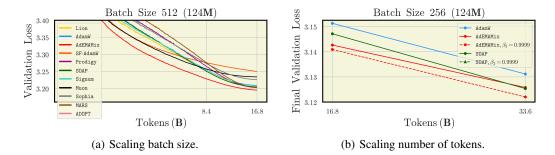


Figure 3: Our results demonstrate that (a): scaling the batch size significantly improves MARS, Signum, Lion and Prodigy making them as good as AdamW even for a long training for $16.8\mathbf{B}$ tokens. Which was not the case in Figure 2(b), where we still observed a significant gap in performance; and (b): indeed, with scaling of the number of iterations, the gap between SOAP and AdEMAMix narrow and, finally, increases. But, on the other hand, with increase of the AdEMAMix β_3 parameter, the performance gap with SOAP reappears.

rate constant proves to be an effective technique for training on shorter horizons. This approach is so effective that methods like Signum and Lion with high weight decay significantly outperform AdamW without weight decay (see Figure 4). Implementation details also warrant attention. Coupled weight decay is still used in some settings, including the PyTorch 63 optimizer implementations. Notably, the popular implementation of Signum becomes ineffective when weight decay is applied. Highlighting this oversight for the community, we contribute by demonstrating our implementation of Signum (Algorithm 6) with decoupled weight decay. The influence of weight decay on model weights is intriguing. As is known, model weights typically grow during training, but weight decay, by modifying the optimized function, significantly reduces the growth of the model's parameter norm. Such ablations of weight decay are also of interest to the community 13 40.

Regarding the ablation of weight decay for optimizers, we again select the best setup for each and conduct a sweep over weight decay values. Our results are presented in Figure 4 and in Figure 21.

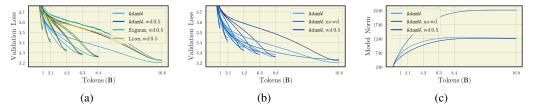


Figure 4: Larger weight decay achieves significantly better results when training on fewer tokens. In (a) we observe that runs of AdamW, Signum, and Lion with the large weight decay of 0.5 consistently outperform the baseline AdamW with weight decay of 0.1 for all training durations except for the last one. Notably, Signum and Lion with large weight decay perform even better than AdamW with the same learning rate. In (b), we also consider a setting without weight decay. We observe that this is suboptimal not only for AdamW, but also for the majority of other optimizers (see Appendix E.2), while the typical weight decay of 0.1 remains the best for large training durations. Importantly, in (c), we ablate the impact od weight decay on the model's ℓ_2 norm.

With our weight decay ablation, we are ready to provide one more insight.

Takeaway 3. The use of weight decay, particularly a large decoupled weight decay term, can significantly impact the final loss value and optimizer behavior. However, for extended training horizons, a moderate, non-zero weight decay proves to be a robust option.

Learning rate sensitivity. Since we tune optimizers at a smaller scale and then extrapolate, we pose the question whether the best learning rate we have found so far transfers to the larger training duration. To verify this, we run 124M model on 16.8B tokens in 256×512 batch size setting, sweeping the learning rate across five typical values: $\{1e^{-4}, 3e^{-4}, 5e^{-4}, 1e^{-3}, 2e^{-3}\}$. The best

learning rate for each method at the moment of hyperparameter tuning on near Chinchilla-optimal 2.1B training duration we report in Appendix D.1 A summary of our results for larger number of tokens is provided in Figure 5 and detailed results of the sweep are presented in Appendix E.2

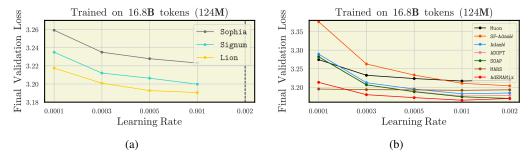


Figure 5: **Optimal learning rate stability across optimizers.** The optimal learning rate determined during tuning on 2.1B tokens remains consistent after a learning rate sweep on 16.8B tokens for most optimizers. In (a), we observe that sign-based methods and similar to them Sophia diverge with increasing learning rate. Interestingly, in (b), SF-AdamW and SOAP demonstrate the best performance with a large learning rate of 0.002. In our work, we further show that it is possible to increase the learning rate even more for such methods.

Warmup ablation. Another important ingredient of the pretraining is learning-rate warmup in the initial phase of training. Recent studies have explored the necessity of warmup in modern deep learning, with some investigating its elimination [39] and others ablating it to improve model performance and stability [92]. We focus on the latter, examining how warmup affects optimizer setup and whether it can significantly enhance performance. For each optimizer's best configuration, we vary warmup across three values: $\{0.27, 1, 4.2\}$ B tokens, which corresponds to $\{2, 8, 32\}$ k iterations. Our choice of the largest warmup value is inspired by [92]. We describe this experiment in Appendix [E.2] Mainly, we observe that Signum and SF-AdamW perform better with a larger warmup of 8k steps when training on 16.8B tokens. We also ablate the claim from [92] that a warmup of 25% of the Chinchilla optimal duration is the best. However, our findings contradict this assertion (see Figure [18]). We show that a moderate values of the warmup, generally, is better, however, different optimizers could prefer different number of warmup steps. As such, SF-AdamW, Sophia, Signum prefer larger warmup, which is clearly depicted in Figure [6].

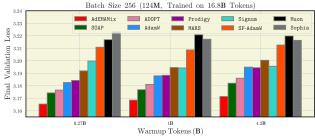


Figure 6: **Warmup ablation.** We report the final validation loss on the FineWeb dataset for 124M model trained on the batch size of 256. We sweep over the batch sizes of $\{1.56\%, 6.25\%, 25\%\}$ of the length of training, which corresponds to $\{2000, 8000, 32000\}$ k iterations, respectively.

Cosine vs WSD. At the outset of our study, we indicated a preference for the cosine scheduler over WSD. In this section, we provide a more detailed ablation of this choice. Having optimally tuned the cosine scheduler for each optimizer, we replicate the setup of [28], which allows us to avoid adjusting additional hyperparameters. Our findings, which demonstrate the superiority of the cosine scheduler across various optimization methods, are presented in Figure [7] and in the Appendix Figures [23] and [24]. These results not only validate our initial preference but also provide insights into the interaction between learning rate schedules and different optimizers in large-scale language model training.

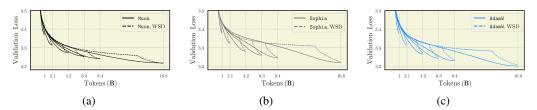


Figure 7: **Comparisons between WSD and cosine scheduler.** Notably, WSD and cosine scheduler behave differently with respect to optimizer. In (a), the Muon optimizer shows a preference for WSD across most training durations. Sophia exhibits an almost perfect match between both schedulers. However, for AdamW, along with the majority of other optimizers studied (see Figure 24), we get a better performance with cosine. We also report a detailed comparison with linear scheduler in Appendix E.2.

4.2 Benchmarking at medium scale: Training Models of 210M Parameters

260

261

263

264 265

266

267

268

269

In this section, we verify if our selected hyperparameters from smaller 124M allow accurate transfer to a slightly larger model. We point out that the most important hyperparameters to be sweeped are learning rate and gradient clipping. Regarding the learning rate, we observe that it only becomes a sensitive choice for sign-based methods, while the optimal hyperparameters for AdamW remain the same.

Results with a batch size of 256×512 . Results provided in the Appendix in Figure 20 are consistent with those obtained training 124M models with large batches.

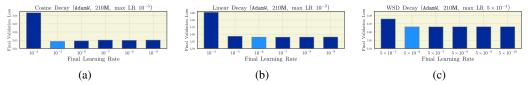


Figure 8: Decaying the learning rate down to $0.01 \times \gamma_{\rm max}$ and beyond, instead of only to 10% We observe a common pattern for different schedulers that decreasing the learning rate to moderate 10^{-2} value is a better choice than decreasing it down to zero. Interestingly, the linear learning rate scheduler for models at a given scale, requires $0.001 \times \gamma_{\rm max}$.

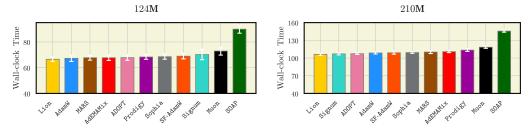


Figure 9: Wall-clock time comparison. After conducting experiments for 124M and 210M models, we are ready to present the wall-time comparison for each methods. For this purposes, we use a single GPU, and run each optimizer for 100 iterations on a small batch size without gradient accumulation and torch.compile. We report the wall-clock time per 100 iterations. We observe that all methods take the roughly the same time or very close time to complete 100 iterations, with the exception of Muon and SOAP. In addition, we point out that SOAP's runtime exhibits a non-linear dependence on the model size, due to its preconditioner matrices operations which are fast only for certain matrices smaller than a predefined size.

4.3 Scaling Up: Benchmarking models of 583M and 720M Parameters

We pick three methods: AdamW, SOAP, and AdEMAMix, and run experiments with a larger model of 583M parameters, and a large batch size of 2M tokens. The goal being to get closer to one of the

settings described in [84]. We train for 6500 and 16000 iterations, corresponding to 13B and 32B tokens respectively.

Comparison between our setting and [84]. We found several key differences between our codebase and [84]: (i) we decay the learning rate to $0.01 \times \gamma$ instead of $0.1 \times \gamma$, with γ being the maximum learning rate, (ii) we use typical weight decay values of e.g. 0.1 instead of smaller values such as 0.01 or 0.0001, (iii) we do not use a z-loss in addition to ours. It has been shown recently that properly decaying the learning rate has an important effect on the optimization [4]. We run an ablation to compare both settings and conclude that removing the z-loss and increasing the weight decay to 0.1 improves the results. Results further improve when the learning rate is decayed more. This ablation is shown in Figure [8]

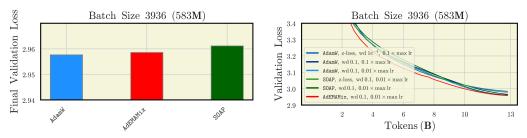


Figure 10: **Results for the** 583**M model.** On the **left**, we show our results when training for 6500 iterations. In this setting, AdamW gives best results, followed by AdEMAMix and then SOAP. This is surprising as it conflicts with findings from [84]. Those results are partly reconciled with the figure on the **right**. And we see that the difference in performance between models trained with and without the z-loss regularizer is quite minor.

5 Extension to MoEs

Setup & Comparison. Besides training dense Llama-like transformers, we also conver a comparison for MoE architectures [78]. Our variant of MoE is based on the Switch-Transformer implementation [20]. We use a classical linear gating with softmax and top-k routing (k=2) and 8 experts. The activation functions remains the same as for the dense base model from Section [3] Such a configuration of the MoE model gives us approximately 520M parameters. We cover additional details in Appendix [E.6] In this setting we train with a batch size of 256×512 for $T \in \{42, 336\}$ k iterations. Again we cover both a Chinchilla-optimal horizon and the beyond. We summarize the results in the following Table [1].

Opt.	42 k	$336\mathbf{k}$
AdEMAMix	22.37	18.47
D-Muon	22.67	18.51
ADOPT	22.70	18.58
AdamW	22.85	18.69
Prodigy	22.82	18.78

Opt.	$42\mathbf{k}$	$336\mathbf{k}$
Lion	23.20	18.87
Signum	23.31	19.09
SF-AdamW	23.34	19.13
Sophia	23.41	19.22
MARS	22.73	19.33

Table 1: Final validation perplexity for MoE training (\downarrow) .

6 Discussion

Our advices on tuning each method. Overall, we validate the already widely used hyperparameters of $\lambda=0.1$ and $T_{\rm warmup}\approx 2{\bf k}$. For Lion—as mentioned in D—we find that the best value for β_1 is consistently 0.99. The mechanism for Lion seems similar to AdEMAMix, one can imagine that Lion could be better with larger β_1 , which would require schedulers. We also pose an interesting observation toward Prodigy: while it may not be so efficient with a super small batch sizes, with scaling of the model size and the batch size it becomes almost as competitive as AdamW. Importantly, Muon and D-Muon performed poorly at a small scale with relatively small batch sizes (32, 256), however, as we see in Figure 1

References

- 1300 [1] Kwangjun Ahn and Byron Xu. Dion: A communication-efficient optimizer for large models, 2025.
- [2] Lukas Balles and Philipp Hennig. Dissecting adam: The sign, magnitude and variance ofstochastic gradients, 2020.
- [3] C. Bekas, E. Kokiopoulou, and Y. Saad. An estimator for the diagonal of a matrix. *Applied Numerical Mathematics*, 57(11):1214–1229, 2007. Numerical Algorithms, Parallelism and Applications (2).
- Shane Bergsma, Nolan Dey, Gurpreet Gosal, Gavia Gray, Daria Soboleva, and Joel Hestness. Straight to zero: Why linearly decaying the learning rate to zero works best for llms, 2025.
- [5] Jeremy Bernstein and Laker Newhouse. Old optimizer, new norm: An anthology, 2024.
- [6] Jeremy Bernstein, Yu-Xiang Wang, Kamyar Azizzadenesheli, and Anima Anandkumar. signsgd: Compressed optimisation for non-convex problems, 2018.
- [7] Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal,
 Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel
 Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M.
 Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz
 Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec
 Radford, Ilya Sutskever, and Dario Amodei. Language models are few-shot learners, 2020.
- [8] David Edwin Carlson, Edo Collins, Ya-Ping Hsieh, Lawrence Carin, and Volkan Cevher.
 Preconditioned spectral descent for deep learning. In *Neural Information Processing Systems*,
 2015.
- [9] Xiangning Chen, Chen Liang, Da Huang, Esteban Real, Kaiyuan Wang, Yao Liu, Hieu Pham,
 Xuanyi Dong, Thang Luong, Cho-Jui Hsieh, Yifeng Lu, and Quoc V. Le. Symbolic discovery
 of optimization algorithms, 2023.
- [10] Savelii Chezhegov, Yaroslav Klyukin, Andrei Semenov, Aleksandr Beznosikov, Alexander
 Gasnikov, Samuel Horváth, Martin Takáč, and Eduard Gorbunov. Gradient clipping improves
 adagrad when the noise is heavy-tailed, 2024.
- 1327 [11] Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul Barham, and Hyung Won. Palm: Scaling language modeling with pathways, 2022.
- George E. Dahl, Frank Schneider, Zachary Nado, Naman Agarwal, Chandramouli Shama Sastry,
 Philipp Hennig, Sourabh Medapati, Runa Eschenhagen, Priya Kasimbeg, Daniel Suo, Juhan
 Bae, Justin Gilmer, Abel L. Peirson, Bilal Khan, Rohan Anil, Mike Rabbat, Shankar Krishnan,
 Daniel Snider, Ehsan Amid, Kongtao Chen, Chris J. Maddison, Rakshith Vasudev, Michal
 Badura, Ankush Garg, and Peter Mattson. Benchmarking Neural Network Training Algorithms,
 2023.
- Francesco D'Angelo, Maksym Andriushchenko, Aditya Varre, and Nicolas Flammarion. Why do we need weight decay in modern deep learning?, 2024.
- 1338 [14] Tri Dao, Daniel Y. Fu, Stefano Ermon, Atri Rudra, and Christopher Ré. Flashattention: Fast and memory-efficient exact attention with io-awareness, 2022.
- 1340 [15] DeepSeek-AI. Deepseek-v3 technical report, 2024.
- 341 [16] Aaron Defazio and Konstantin Mishchenko. Learning-rate-free learning by d-adaptation, 2023.
- [17] Aaron Defazio, Xingyu Alice Yang, Harsh Mehta, Konstantin Mishchenko, Ahmed Khaled, and
 Ashok Cutkosky. The road less scheduled, 2024.

- 18 Nolan Dey, Quentin Anthony, and Joel Hestness. The practitioner's guide to the maximal update parameterization. https://www.cerebras.ai/blog/the-practitioners-guide-to-the-maximal-update-parameterization. September 2024.
- [19] John Duchi, Elad Hazan, and Yoram Singer. Adaptive subgradient methods for online learning
 and stochastic optimization. *Journal of Machine Learning Research*, 12(61):2121–2159, 2011.
- [20] William Fedus, Barret Zoph, and Noam Shazeer. Switch transformers: Scaling to trillion
 parameter models with simple and efficient sparsity, 2022.
- Leo Gao, Stella Biderman, Sid Black, Laurence Golding, Travis Hoppe, Charles Foster, Jason
 Phang, Horace He, Anish Thite, Noa Nabeshima, Shawn Presser, and Connor Leahy. The Pile:
 An 800gb dataset of diverse text for language modeling. arXiv preprint arXiv:2101.00027,
 2020.
- 356 [22] Google Gemini Team. Gemini: A family of highly capable multimodal models, 2024.
- 357 [23] Alex Graves. Generating sequences with recurrent neural networks, 2014.
- Vineet Gupta, Tomer Koren, and Yoram Singer. Shampoo: Preconditioned stochastic tensor optimization, 2018.
- [25] Nicholas J. Higham. Functions of Matrices. Society for Industrial and Applied Mathematics,
 2008.
- [26] Jordan Hoffmann, Sebastian Borgeaud, Arthur Mensch, Elena Buchatskaya, Trevor Cai, Eliza
 Rutherford, Diego de Las Casas, Lisa Anne Hendricks, Johannes Welbl, Aidan Clark, Tom
 Hennigan, Eric Noland, Katie Millican, George van den Driessche, Bogdan Damoc, Aurelia
 Guy, Simon Osindero, Karen Simonyan, Erich Elsen, Jack W. Rae, Oriol Vinyals, and Laurent
 Sifre. Training compute-optimal large language models, 2022.
- Shengding Hu, Yuge Tu, Xu Han, Chaoqun He, Ganqu Cui, Xiang Long, Zhi Zheng, Yewei
 Fang, Yuxiang Huang, Weilin Zhao, Xinrong Zhang, Zheng Leng Thai, Kaihuo Zhang, Chongyi
 Wang, Yuan Yao, Chenyang Zhao, Jie Zhou, Jie Cai, Zhongwu Zhai, Ning Ding, Chao Jia,
 Guoyang Zeng, Dahai Li, Zhiyuan Liu, and Maosong Sun. Minicpm: Unveiling the potential of
 small language models with scalable training strategies, 2024.
- ³⁷² [28] Alexander Hägele, Elie Bakouch, Atli Kosson, Loubna Ben Allal, Leandro Von Werra, and Martin Jaggi. Scaling laws and compute-optimal training beyond fixed training durations, 2024.
- 274 [29] Pavel Izmailov, Dmitrii Podoprikhin, Timur Garipov, Dmitry Vetrov, and Andrew Gordon Wilson. Averaging weights leads to wider optima and better generalization, 2019.
- [30] Sami Jaghouar, Jack Min Ong, Manveer Basra, Fares Obeid, Jannik Straube, Michael Keiblinger,
 Elie Bakouch, Lucas Atkins, Maziyar Panahi, Charles Goddard, Max Ryabinin, and Johannes
 Hagemann. Intellect-1 technical report, 2024.
- 379 [31] Albert Q. Jiang, Alexandre Sablayrolles, Antoine Roux, Arthur Mensch, Blanche Savary, Chris
 380 Bamford, Devendra Singh Chaplot, Diego de las Casas, Emma Bou Hanna, Florian Bressand,
 381 Gianna Lengyel, Guillaume Bour, Guillaume Lample, Lélio Renard Lavaud, Lucile Saulnier,
 382 Marie-Anne Lachaux, Pierre Stock, Sandeep Subramanian, Sophia Yang, Szymon Antoniak,
 383 Teven Le Scao, Théophile Gervet, Thibaut Lavril, Thomas Wang, Timothée Lacroix, and
 384 William El Sayed. Mixtral of experts, 2024.
- [32] Keller Jordan, Jeremy Bernstein, Brendan Rappazzo, @fernbear.bsky.social, Boza Vlado,
 You Jiacheng, Franz Cesista, Braden Koszarsky, and @Grad62304977. modded-nanogpt:
 Speedrunning the nanogpt baseline, 2024.
- [33] Keller Jordan, Yuchen Jin, Vlado Boza, You Jiacheng, Franz Cecista, Laker Newhouse, and
 Jeremy Bernstein. Muon: An optimizer for hidden layers in neural networks, 2024.
- [34] Jean Kaddour, Oscar Key, Piotr Nawrot, Pasquale Minervini, and Matt J. Kusner. No train no
 gain: Revisiting efficient training algorithms for transformer-based language models, 2023.

- [35] Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B. Brown, Benjamin Chess, Rewon Child,
 Scott Gray, Alec Radford, Jeffrey Wu, and Dario Amodei. Scaling laws for neural language
 models, 2020.
- Sai Praneeth Karimireddy, Quentin Rebjock, Sebastian Stich, and Martin Jaggi. Error feed back fixes signSGD and other gradient compression schemes. In *ICML 2019 International Conference on Machine Learning*, pages 3252–3261. PMLR, 2019.
- 398 [37] Andrej Karpathy. NanoGPT. https://github.com/karpathy/nanoGPT, 2022.
- 399 [38] Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization, 2017.
- 400 [39] Atli Kosson, Bettina Messmer, and Martin Jaggi. Analyzing & reducing the need for learning rate warmup in gpt training, 2024.
- 402 [40] Atli Kosson, Bettina Messmer, and Martin Jaggi. Rotational equilibrium: How weight decay balances learning across neural networks, 2024.
- [41] Frederik Kunstner. Why do machine learning optimizers that work, work? PhD thesis, University of British Columbia, 2024.
- 406 [42] Frederik Kunstner, Jacques Chen, Jonathan Wilder Lavington, and Mark Schmidt. Noise is not
 407 the main factor behind the gap between sgd and adam on transformers, but sign descent might
 408 be, 2023.
- [43] Frederik Kunstner, Robin Yadav, Alan Milligan, Mark Schmidt, and Alberto Bietti. Heavy-tailed class imbalance and why adam outperforms gradient descent on language models, 2024.
- [44] Jeffrey Li, Alex Fang, Georgios Smyrnis, Maor Ivgi, Matt Jordan, Samir Yitzhak Gadre, Hritik
 Bansal, Etash Guha, Sedrick Scott Keh, Kushal Arora, et al. Datacomp-lm: In search of the next
 generation of training sets for language models. Advances in Neural Information Processing
 Systems, 37:14200–14282, 2024.
- [45] Shen Li, Yanli Zhao, Rohan Varma, Omkar Salpekar, Pieter Noordhuis, Teng Li, Adam Paszke,
 Jeff Smith, Brian Vaughan, Pritam Damania, and Soumith Chintala. Pytorch distributed:
 Experiences on accelerating data parallel training, 2020.
- 418 [46] Hong Liu, Zhiyuan Li, David Hall, Percy Liang, and Tengyu Ma. Sophia: A scalable stochastic second-order optimizer for language model pre-training, 2024.
- [47] Jingyuan Liu, Jianlin Su, Xingcheng Yao, Zhejun Jiang, Guokun Lai, Yulun Du, Yidao Qin,
 Weixin Xu, Enzhe Lu, Junjie Yan, Yanru Chen, Huabin Zheng, Yibo Liu, Shaowei Liu, Bohong
 Yin, Weiran He, Han Zhu, Yuzhi Wang, Jianzhou Wang, Mengnan Dong, Zheng Zhang,
 Yongsheng Kang, Hao Zhang, Xinran Xu, Yutao Zhang, Yuxin Wu, Xinyu Zhou, and Zhilin
 Yang. Muon is scalable for llm training, 2025.
- 425 [48] AI @ Meta Llama Team. The llama 3 herd of models, 2024.
- 426 [49] Ilya Loshchilov and Frank Hutter. Sgdr: Stochastic gradient descent with warm restarts, 2017.
- 427 [50] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization, 2019.
- 428 [51] James Martens. New insights and perspectives on the natural gradient method, 2020.
- [52] James Martens and Roger Grosse. Optimizing neural networks with kronecker-factored approximate curvature. In *International conference on machine learning*, pages 2408–2417. PMLR, 2015.
- Faulius Micikevicius, Sharan Narang, Jonah Alben, Gregory Diamos, Erich Elsen, David Garcia,
 Boris Ginsburg, Michael Houston, Oleksii Kuchaiev, Ganesh Venkatesh, and Hao Wu. Mixed
 precision training, 2018.
- 435 [54] Konstantin Mishchenko and Aaron Defazio. Prodigy: An expeditiously adaptive parameter-free learner, 2024.

- 437 [55] A.S. Nemirovskii and Yu.E. Nesterov. Optimal methods of smooth convex minimization. *USSR Computational Mathematics and Mathematical Physics*, 25(2):21–30, 1985.
- 439 [56] Yu. Nesterov and V. Shikhman. Quasi-monotone Subgradient Methods for Nonsmooth Convex Minimization. *Journal of Optimization Theory and Applications*, 165(3):917–940, June 2015.
- 441 [57] Yurii Nesterov. A method for unconstrained convex minimization problem with the rate of convergence $o(1/k^2)$, 1983.
- 443 [58] Team OLMo. 2 olmo 2 furious, 2024.
- 444 [59] OpenAI. Gpt-4 technical report, 2024.
- 445 [60] Francesco Orabona. Neural networks (maybe) evolved to make adam the best optimizer, 2020.
- [61] Francesco Orabona and Dávid Pál. Open problem: Parameter-free and scale-free online algorithms. In Vitaly Feldman, Alexander Rakhlin, and Ohad Shamir, editors, 29th Annual
 Conference on Learning Theory, volume 49 of Proceedings of Machine Learning Research, pages 1659–1664, Columbia University, New York, New York, USA, 23–26 Jun 2016. PMLR.
- 450 [62] Matteo Pagliardini, Pierre Ablin, and David Grangier. The ademamix optimizer: Better, faster, older, 2024.
- [63] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan,
 Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas
 Köpf, Edward Yang, Zach DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy,
 Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. Pytorch: An imperative style,
 high-performance deep learning library, 2019.
- Guilherme Penedo, Hynek Kydlíček, Loubna Ben allal, Anton Lozhkov, Margaret Mitchell,
 Colin Raffel, Leandro Von Werra, and Thomas Wolf. The fineweb datasets: Decanting the web
 for the finest text data at scale, 2024.
- [65] Guilherme Penedo, Hynek Kydlíček, Vinko Sabolčec, Bettina Messmer, Negar Foroutan, Martin
 Jaggi, Leandro von Werra, and Thomas Wolf. Fineweb2: A sparkling update with 1000s of
 languages, December 2024.
- 463 [66] Thomas Pethick, Wanyun Xie, Kimon Antonakopoulos, Zhenyu Zhu, Antonio Silveti-Falls, and Volkan Cevher. Training deep learning models with norm-constrained lmos, 2025.
- 465 [67] Boris Polyak. Some methods of speeding up the convergence of iteration methods. *Ussr Computational Mathematics and Mathematical Physics*, 4:1–17, 1964.
- [68] Boris Polyak. New method of stochastic approximation type. *Automation and Remote Control*,1990, 01 1990.
- Tomer Porian, Mitchell Wortsman, Jenia Jitsev, Ludwig Schmidt, and Yair Carmon. Resolving discrepancies in compute-optimal scaling of language models, 2024.
- 471 [70] Alec Radford and Karthik Narasimhan. Improving language understanding by generative pre-training. 2018.
- 473 [71] Alec Radford, Jeff Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. Language models are unsupervised multitask learners. *OpenAI*, 2019.
- 475 [72] Herbert Robbins and Sutton Monro. A Stochastic Approximation Method. *The Annals of Mathematical Statistics*, 22(3):400 407, 1951.
- 477 [73] David Ruppert. Efficient estimations from a slowly convergent robbins-monro process. 1988.
- 478 [74] Nikhil Sardana, Jacob Portes, Sasha Doubov, and Jonathan Frankle. Beyond chinchilla-optimal:
 479 Accounting for inference in language model scaling laws, 2024.
- Fabian Schaipp, Alexander Hägele, Adrien Taylor, Umut Simsekli, and Francis Bach. The surprising agreement between convex optimization theory and learning-rate scheduling for large model training, 2025.

- 483 [76] Robin M. Schmidt, Frank Schneider, and Philipp Hennig. Descending through a crowded valley
 484 benchmarking deep learning optimizers, 2021.
- ⁴⁸⁵ [77] Noam Shazeer. Glu variants improve transformer, 2020.
- [78] Noam Shazeer, Azalia Mirhoseini, Krzysztof Maziarz, Andy Davis, Quoc Le, Geoffrey Hinton,
 and Jeff Dean. Outrageously large neural networks: The sparsely-gated mixture-of-experts
 layer, 2017.
- Yikang Shen, Matthew Stallone, Mayank Mishra, Gaoyuan Zhang, Shawn Tan, Aditya Prasad,
 Adriana Meza Soria, David D. Cox, and Rameswar Panda. Power scheduler: A batch size and
 token number agnostic learning rate scheduler, 2024.
- 492 [80] Charlie Snell, Jaehoon Lee, Kelvin Xu, and Aviral Kumar. Scaling Ilm test-time compute optimally can be more effective than scaling model parameters, 2024.
- [81] Jianlin Su, Yu Lu, Shengfeng Pan, Ahmed Murtadha, Bo Wen, and Yunfeng Liu. Roformer:
 Enhanced transformer with rotary position embedding, 2023.
- [82] Ilya Sutskever, James Martens, George Dahl, and Geoffrey Hinton. On the importance of initialization and momentum in deep learning. In Sanjoy Dasgupta and David McAllester, editors, *Proceedings of the 30th International Conference on Machine Learning*, Proceedings of Machine Learning Research, pages 1139–1147, Atlanta, Georgia, USA, 17–19 Jun 2013.
 PMLR.
- 501 [83] Shohei Taniguchi, Keno Harada, Gouki Minegishi, Yuta Oshima, Seong Cheol Jeong, Go Na-502 gahara, Tomoshi Iiyama, Masahiro Suzuki, Yusuke Iwasawa, and Yutaka Matsuo. Adopt: 503 Modified adam can converge with any β_2 with the optimal rate, 2024.
- [84] Nikhil Vyas, Depen Morwani, Rosie Zhao, Itai Shapira, David Brandfonbrener, Lucas Janson,
 and Sham Kakade. Soap: Improving and stabilizing shampoo using adam, 2024.
- [85] Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed Chi,
 Quoc Le, and Denny Zhou. Chain-of-thought prompting elicits reasoning in large language
 models, 2023.
- [86] Aiyuan Yang, Bin Xiao, Bingning Wang, Borong Zhang, Ce Bian, Chao Yin, Chenxu Lv,
 Da Pan, Dian Wang, Dong Yan, et al. Baichuan 2: Open large-scale language models. arXiv
 preprint arXiv:2309.10305, 2023.
- [87] Greg Yang, Edward J. Hu, Igor Babuschkin, Szymon Sidor, Xiaodong Liu, David Farhi, Nick
 Ryder, Jakub Pachocki, Weizhu Chen, and Jianfeng Gao. Tensor programs v: Tuning large
 neural networks via zero-shot hyperparameter transfer, 2022.
- Huizhuo Yuan, Yifeng Liu, Shuang Wu, Xun Zhou, and Quanquan Gu. Mars: Unleashing the power of variance reduction for training large models, 2024.
- [89] Manzil Zaheer, Sashank Reddi, Devendra Sachan, Satyen Kale, and Sanjiv Kumar. Adaptive
 methods for nonconvex optimization. In S. Bengio, H. Wallach, H. Larochelle, K. Grauman,
 N. Cesa-Bianchi, and R. Garnett, editors, Advances in Neural Information Processing Systems,
 volume 31. Curran Associates, Inc., 2018.
- [90] Xiaohua Zhai, Alexander Kolesnikov, Neil Houlsby, and Lucas Beyer. Scaling vision transformers, 2022.
- 523 [91] Biao Zhang and Rico Sennrich. Root mean square layer normalization, 2019.
- [92] Hanlin Zhang, Depen Morwani, Nikhil Vyas, Jingfeng Wu, Difan Zou, Udaya Ghai, Dean
 Foster, and Sham Kakade. How does critical batch size scale in pre-training?, 2024.
- [93] Jingzhao Zhang, Sai Praneeth Karimireddy, Andreas Veit, Seungyeon Kim, Sashank Reddi,
 Sanjiv Kumar, and Suvrit Sra. Why are adaptive methods good for attention models? *Advances in Neural Information Processing Systems*, 33:15383–15393, 2020.

- [94] Jingzhao Zhang, Sai Praneeth Karimireddy, Andreas Veit, Seungyeon Kim, Sashank J Reddi,
 Sanjiv Kumar, and Suvrit Sra. Why are adaptive methods good for attention models?, 2020.
- [95] Yushun Zhang, Congliang Chen, Tian Ding, Ziniu Li, Ruoyu Sun, and Zhi-Quan Luo. Why
 transformers need adam: A hessian perspective, 2024.
- Fig. 1963 Rosie Zhao, Depen Morwani, David Brandfonbrener, Nikhil Vyas, and Sham Kakade. Deconstructing what makes a good optimizer for language models. *ICLR*, 2025.

NeurIPS Paper Checklist

1. Claims

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

Answer: [Yes]

Justification: the contribution of this paper is described accurately in the abstract and introduction.

Guidelines:

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

2. Limitations

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

Answer: [Yes]

Justification: we discuss a limitations and mention experiments we have not tried to run

- Guidelines:

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

3. Theory assumptions and proofs

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

Answer: [NA]

Justification: this is not a theoretical work 587 Guidelines: 588 589 590

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

4. Experimental result reproducibility

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

Answer: [Yes]

591

592

593

594

595

596

597

598

599

600

601

602

603

604

605

606

607

608

609

610

611

612

613

614

615

616

617

618

619

620

621 622

623

625

626

627

628

629

630

631

632

633

634

635

636

637

638

639

640

Justification: we open-source our code

Guidelines:

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

5. Open access to data and code

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

641	Answer: [Yes]
642	Justification: we provide our code and the datasets are mentioned clearly
643	Guidelines:
644	• The answer NA means that paper does not include experiments requiring code.
645	• Please see the NeurIPS code and data submission guidelines (https://nips.cc/
646	public/guides/CodeSubmissionPolicy) for more details.
647	• While we encourage the release of code and data, we understand that this might not be

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

6. Experimental setting/details

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

Answer: [Yes]

Justification: we specify all of the important hyperparameters as well as hyperparameter tuning.

Guidelines:

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

7. Experiment statistical significance

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

Answer: [No]

Justification: in our large-scale experiments we could not affort so. and we are running all of the experiment with the same seed for generation data splits, etc.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The authors should answer "Yes" if the results are accompanied by error bars, confidence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.
- The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, random drawing of some parameter, or overall run with given experimental conditions).
- The method for calculating the error bars should be explained (closed form formula, call to a library function, bootstrap, etc.)
- The assumptions made should be given (e.g., Normally distributed errors).

- It should be clear whether the error bar is the standard deviation or the standard error
 of the mean.
 - It is OK to report 1-sigma error bars, but one should state it. The authors should preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis of Normality of errors is not verified.
 - For asymmetric distributions, the authors should be careful not to show in tables or figures symmetric error bars that would yield results that are out of range (e.g. negative error rates).
 - If error bars are reported in tables or plots, The authors should explain in the text how
 they were calculated and reference the corresponding figures or tables in the text.

8. Experiments compute resources

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

Answer: [Yes]

Justification: we provide this in Appendix

Guidelines:

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

9. Code of ethics

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

Answer: [Yes]

Justification: this paper is consistent with NeurIPS Code of Ethics.

Guidelines:

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

10. Broader impacts

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

Answer: [NA]

Justification: there is no societal impact of the work performed.

Guidelines:

- The answer NA means that there is no societal impact of the work performed.
- If the authors answer NA or No, they should explain why their work has no societal impact or why the paper does not address societal impact.
- Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.

- The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
- The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
- If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

11. Safeguards

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

Answer: [NA]

Justification: the paper poses no such risks.

Guidelines:

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

12. Licenses for existing assets

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

Answer: [Yes]

Justification: we cite them and respect, see Sections 1 and 2

Guidelines:

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

 If this information is not available online, the authors are encouraged to reach out to the asset's creators.

13. New assets

794

795

796

797

798

799

800

801

802

803

804

805

806

807

808

809

810

811

812

813

814

815

816

818

819

820 821

822

823

824

825

826

827

828

830

831

832

833

834

835

836

837 838

839

840

841

842

843

844

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

Answer: [NA]

Justification: the paper does not propose new assets.

Guidelines:

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

14. Crowdsourcing and research with human subjects

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

Answer: [NA]

Justification: our work does not include research with human subjects.

Guidelines:

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

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

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

Answer: [NA]

Justification: our work does not include research with human subjects.

Guidelines:

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

16. Declaration of LLM usage

Question: Does the paper describe the usage of LLMs if it is an important, original, or 845 non-standard component of the core methods in this research? Note that if the LLM is used 846 only for writing, editing, or formatting purposes and does not impact the core methodology, 847 scientific rigorousness, or originality of the research, declaration is not required. 848 Answer: [NA] 849 Justification: the core development of our work does not involve LLMs. 850 Guidelines: 851 • The answer NA means that the core method development in this research does not 852 involve LLMs as any important, original, or non-standard components. 853

854

855

• Please refer to our LLM policy (https://neurips.cc/Conferences/2025/LLM)

for what should or should not be described.