

# Fantastic Pretraining Optimizers and Where to Find Them

Kaiyue Wen  
Stanford University  
kaiyuew@stanford.edu

David Hall  
Stanford University  
dlwh@cs.stanford.edu

Tengyu Ma  
Stanford University  
tengyuma@stanford.edu

Percy Liang  
Stanford University  
плианг@cs.stanford.edu

September 3, 2025

## Abstract

AdamW has long been the dominant optimizer in language model pretraining, despite numerous claims that alternative optimizers offer  $1.4\times$  to  $2\times$  speedup. We posit that two methodological shortcomings have obscured fair comparisons and hindered practical adoption: (i) unequal hyperparameter tuning and (ii) limited or misleading evaluation setups. To address these two issues, we conduct a systematic study of ten deep learning optimizers across four model scales (0.1B-1.2B parameters) and data-to-model ratios ( $1-8\times$  the Chinchilla optimum). We find that fair and informative comparisons require rigorous hyperparameter tuning and evaluations across a range of model scales and data-to-model ratios, performed at the end of training. First, optimal hyperparameters for one optimizer may be suboptimal for another, making blind hyperparameter transfer unfair. Second, the actual speedup of many proposed optimizers over well-tuned baselines is lower than claimed and decreases with model size to only  $1.1\times$  for 1.2B parameter models. Thirdly, comparing intermediate checkpoints before reaching the target training budgets can be misleading, as rankings between two optimizers can flip during training due to learning rate decay. Through our thorough investigation, we find that all the fastest optimizers such as Muon and Soap, use matrices as preconditioners — multiplying gradients with matrices rather than entry-wise scalars. However, the speedup of matrix-based optimizers is inversely proportional to model scale, decreasing from  $1.4\times$  over AdamW for 0.1B parameter models to merely  $1.1\times$  for 1.2B parameter models.

## 1 Introduction

Pretraining has been the most computationally expensive component in the training pipeline for large language models, accounting for over 95% of the cost in DeepSeek V3 DeepSeek-AI et al. [2025b], and the additional RL training cost in DeepSeek R1 DeepSeek-AI et al. [2025a] is also comparatively much smaller. Until recently, AdamW has been the standard optimizer. Recent studies have introduced novel optimizers that claim to accelerate pretraining by  $1.4\times$  to  $2\times$  compared to AdamW Liu et al. [2024a], Vyas et al. [2025], Liu et al. [2025a], Yuan et al. [2025], Liang et al. [2025], Wang et al. [2025], Liu et al. [2025c], Pethick et al. [2025], Ma et al. [2025], yet these optimizers have not yet been widely adopted in real-world pretraining DeepSeek-AI et al. [2024], Yang et al. [2025], Grattafiori et al. [2024], with the notable exception of Kimi K2 Team et al. [2025], which uses the Muon-clip optimizer (a variant of Muon Jordan et al. [2024]).

We pinpoint two problems in optimizer evaluation in the evaluation process of these optimizers that both undermine confidence in new methods and limit their practical adoption. First, some baselines suffer from improperly tuned hyperparameters. Second, many experiments are confined to smaller-scale settings, leaving unanswered questions about how these optimizers perform in broader, more realistic scenarios.

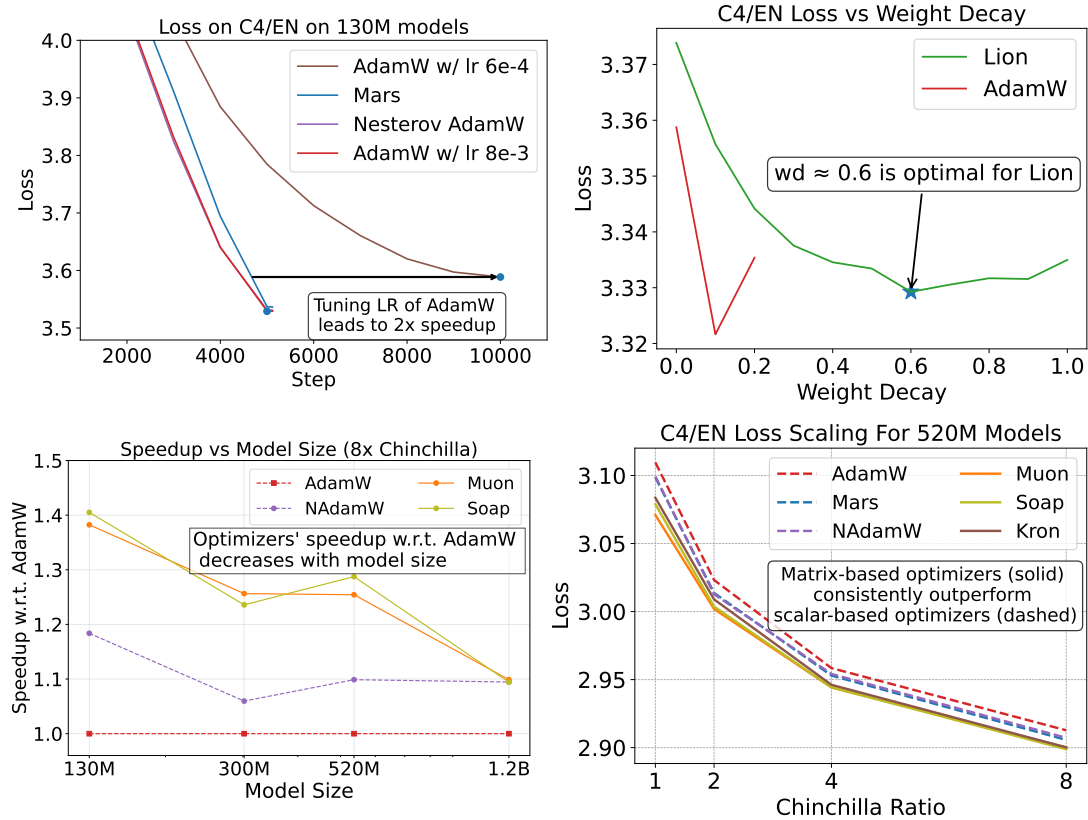


Figure 1: **Top Left:** The commonly used AdamW baseline for optimizer design is under-tuned. Up to a  $2\times$  speedup is achievable by tuning a single hyperparameter (learning rate) in the GPT-3 recipe Brown et al. [2020] for a 100M model (adopted in Liu et al. [2024a], Wen et al. [2024], Yuan et al. [2025], Liang et al. [2025], Wang et al. [2025]), highlighting the importance of proper hyperparameter optimization. **Top Right:** Fixing hyperparameters across optimizers does not guarantee fair comparison. Shared hyperparameters such as learning rate and weight decay are commonly set to a constant in previous studies. However, even conceptually similar optimizers may correspond to very different optimal hyperparameters. **Bottom Left:** Speedup decays with model size. While some optimizers show high ( $1.3\text{--}1.4\times$ ) speedup over AdamW on models under 1B parameters, the speedup decays with model size to only  $1.1\times$  for 1.2B parameters. **Bottom Right:** Matrix-based optimizers consistently outperform scalar-based optimizers. The loss curves for three scalar-based optimizers (AdamW, Nesterov AdamW, Mars) and three matrix-based optimizers (Kron, Soap, Muon) trained with different Chinchilla ratios of data are shown. Matrix-based optimizers achieve a consistent speedup over scalar-based optimizers. Furthermore, the three matrix-based optimizers converge to a similar matrix loss in an overtrained setting.

To address these issues, we benchmark eleven optimizers, including AdamW, in a rigorously controlled setup, focusing on two key questions:

1. **How to ensure hyperparameter optimality?** Previous works typically rely on manual hyperparameter selection and keep common hyperparameters such as learning rate and weight decay fixed across optimizers (e.g. Liu et al. [2025a]). This tuning process may result in a weak baseline, as the hyperparameters chosen may favor the proposed optimizers rather than for AdamW. We validate this concern by showing that tuning just one hyperparameter in the widely adopted GPT-3 recipe (introduced in Brown et al. [2020] and used in Liu et al. [2024a], Wen et al. [2024], Yuan et al. [2025], Liu et al. [2025d,c], Liang et al. [2025], Wang et al. [2025]) can yield a  $2\times$  speedup for pretraining (Figure 1, top left). To address, this we perform coordinate descent in the hyperparameter space on  $N$  hyperparameters for all eleven optimizers, iterating until convergence across for each of six different settings for models with up to 0.5B parameters to ensure that the hyperparameters are near-optimal for each setting.

2. **How do speedups differ across different scaling regimes?** Large language models are trained in different regimes: depending on the model size, the data-to-model ratio (the number of tokens over the number of parameters) can range from 1 to more than 50 times the Chinchilla optimal (about 20 based on Hoffmann et al. [2022]). Typical optimizer experiments in  $1\times$  Chinchilla optimal regimes raise concerns about a new optimizer’s effectiveness in high data-to-model ratio. To address this, we benchmark the optimizers in four distinct data-to-model ratios ( $1\times$ ,  $2\times$ ,  $4\times$  and  $8\times$  the Chinchilla optimal regime) and scale up to 1.2B parameter models (following Everett et al. [2024], Zhang et al. [2025a]).

Concretely, we employ the Llama 2 architecture Touvron et al. [2023b], Grattafiori et al. [2024] (ranging from 0.1B to 1.2B parameters) and a data mixture similar to OLMo 2’s OLMo et al. [2025]. We focus on the final validation loss on the C4-EN mixture as a known proxy for downstream performance Bhagia et al. [2024], while also tracking exact downstream performance on various benchmarks. As previous works Vyas et al. [2025], Liu et al. [2025a] show that the step-wise computation overhead of matrix-based optimizers can be reduced to under 10% through proper implementation, we primarily compare algorithms by the number of tokens needed to reach a given loss.

Our empirical results show the necessity of careful hyperparameter tuning and end-of-training evaluations across a range of model scales and data-to-model ratios: we

1. **Hyperparameter transfer between optimizers is non-trivial.** Even similar optimizers may need very different hyperparameters (e.g., Lion’s optimal weight decay  $\approx 0.6$  vs. AdamW’s  $\approx 0.1$ , Figure 1, top right), so fixing hyperparameters across optimizers can lead to unfair comparisons.
2. **The speedup of new optimizers is lower than claimed and diminishes with model size.** Many reported speedups of  $2\times$  simply reflect a weak baseline. Against our well-tuned AdamW baseline, the speedup of alternative optimizers does not exceed  $1.4\times$  (Figure 3). Furthermore, while new optimizers such as Muon and Soap show  $1.3\times$  speedups for small models (0.1B), the speedups diminish to around  $1.1\times$  for 1.2B parameter models at  $8\times$  Chinchilla Ratio (Figure 1, bottom left), a regime that is not tested in previous works studying the scaling law of these optimizers <sup>1</sup>.
3. **Early-stage loss curves can mislead significantly.** During learning rate decay, loss curves of different optimizers may cross multiple times (Figure 5), so judging optimizers using intermediate checkpoints may result in a different ranking than comparing models at the target training budget.

Our benchmarking also reveals new insights about optimizer design:

1. **Matrix-based optimizers consistently outperform scalar-based optimizers for small models.** *Scalar-based optimizers* (e.g., AdamW, Lion, Mars, etc.) update each parameter individually using scalar operations. After proper tuning, all scalar-based optimizers achieve similar optimization speeds to AdamW, with an average speedup ratio of less than  $1.2\times$ . *Matrix-based optimizers* (e.g., Kron, Muon, Soap, etc.) leverage the inherent matrix structure of neural network parameters and precondition gradients using matrix multiplication. Despite their diverse update rules, matrix-based optimizers all deliver approximately a  $1.3\times$  speedup over AdamW (Figure 1, bottom right) for models under 520M parameters.
2. **Optimal choice of optimizer shifts depends on data-to-model ratios.** A winner in the  $1\times$  Chinchilla regime may be suboptimal when data-to-model ratio increases. For example, while Muon is consistently the best optimizer in smaller Chinchilla ratio regimes, it is outperformed by Kron and Soap when the data-to-model ratio increases to  $8\times$  or larger (Figures 3 and 4).

## 2 Related Works

**Optimizers for Deep Learning.** A long line of work has studied optimization for deep learning, incorporating insights from classical optimization literature Robbins and Monro [1951], Nesterov [1983], Duchi et al. [2011] and domain knowledge about deep neural networks Sutskever et al. [2013]. (i) Early insights motivated the

<sup>1</sup>The original Soap paper Vyas et al. [2025] investigate model sizes up to 0.6B parameters. and Kimi’s paper on Muon Liu et al. [2025a] only considers  $1\times$  Chinchilla regime.

rise of optimizers that use adaptive learning rates based on second-order momentum Tieleman [2012], Zeiler [2012], Kingma and Ba [2017]. Adam Kingma and Ba [2017] later became the default baseline for optimizers with adaptive learning rates. Since then, improvements over Adam and SGD have been proposed, with notable examples including Nesterov Adam Dozat [2016] and AdamW with decoupled weight decay Loshchilov and Hutter [2019]. Other improvements include addressing the convergence of Adam on convex loss Reddi et al. [2018], Zaheer et al. [2018], Taniguchi et al. [2024], considering interpolation between Adam and SGD to improve generalization Luo et al. [2019], Xie et al. [2022], performing further variance reductions on optimizer updates Liu et al. [2021], Zhang et al. [2019], Yuan et al. [2025], Xie et al. [2024], Pagliardini et al. [2024], Zhuang et al. [2020], incorporating momentum on weights Ivgi et al. [2023], Defazio et al. [2024b], allowing easier hyperparameter tuning Defazio and Mishchenko [2023], Mishchenko and Defazio [2024], Defazio et al. [2024a], reducing memory usage by incorporating the structure of neural networks Shazeer and Stern [2018], Zhang et al. [2025b], Zhu et al. [2025], Luo et al. [2023], Modoranu et al. [2024], Zhao et al. [2024], and modifying the algorithm to allow larger batch sizes You et al. [2017, 2020]. (ii) Starting from Preconditioned SGD Li [2018a] and Shampoo Gupta et al. [2018], another line of optimizer design Morwani et al. [2024], Eschenhagen et al. [2023], Martens and Grosse [2020], Li [2018b], Eschenhagen et al. [2025] began to incorporate matrix preconditioners instead of simple scalar preconditioners. Techniques including learning rate grafting Agarwal et al. [2020], blocking, and distributed methodology Anil et al. [2021] have since been proposed. These matrix-based approaches later led to the theory of modular duality in deep learning optimization Bernstein and Newhouse [2024a,b], Large et al. [2024] and new optimizers including Muon Jordan et al. [2024] and Scion Pethick et al. [2025]. (iii) Motivated by Newton’s algorithm, there has been a line of work that tries to incorporate Hessian information Becker and Cun [1989], Yao et al. [2021], Schaul et al. [2013], Yao et al. [2021]. (iv) Symbolic discovery of optimizers Chen et al. [2023] has discovered a memory-efficient SignGD Bernstein et al. [2018] variant called Lion, which claims to outperform Adam on a wide range of tasks.

**Optimization for Pretraining.** Since Brown et al. [2020], the cost of pretraining has increased dramatically. One of the challenges is how to choose hyperparameters with minimal cost. A pivotal line of work in this direction is the tensor program series that allows for extrapolating some hyperparameters across scales Yang [2020b,c,a, 2021], Yang and Littwin [2023], Yang et al. [2022, 2024]. Empirical results suggest that fitting a power law to scale hyperparameters is now common practice for large models Li et al. [2025a], Everett et al. [2024], Liu et al. [2025a], DeepSeek-AI et al. [2024], Zhang et al. [2025a]. We have also incorporated this approach in our paper. Another popular line of research is designing better optimizers specifically for pretraining. These optimizers are our main objects of study. An (incomplete) list of optimizers and their claimed speedups over AdamW includes Sophia Liu et al. [2024a] ( $2\times$ ), Soap Vyas et al. [2025] ( $1.4\times$ ), Muon Jordan et al. [2024], Liu et al. [2025a] ( $2\times$ ), MARS Yuan et al. [2025] ( $2\times$ ), Cautious AdamW Liang et al. [2025] ( $2\times$ ), Block-wise Learning Rate Adam Wang et al. [2025] ( $2\times$ ), FOCUS Liu et al. [2025c] ( $2\times$ ), SWAN Ma et al. [2025] ( $2\times$ ), DION Ahn et al. [2025] ( $3\times$ ), and SPlus Frans et al. [2025] ( $2\times$ ). We present a comparison of setups with these prior works in Appendix G.

**Re-evaluation Methodology.** Our work is a rigorous evaluation of optimizers for pretraining. Rigorous evaluation has been an important part of deep learning research to clarify the current status of research and move the community forward. Jiang et al. [2019] critically examines metrics for predicting LLMs’ generalization capability and has facilitated research on understanding loss landscape sharpness and generalization and new optimizers such as SAM Foret et al. [2021]. Schmidt et al. [2021] re-evaluated optimizers at that time and showed that (i) which optimizer is optimal is problem-specific, and (ii) rigorous hyperparameter tuning is required and unequal tuning can account for most of the claimed speedup. Unlike Schmidt et al. [2021], we evaluate the LLM pretraining task and include modern optimizers they did not test. However, we arrive at a similar conclusion: rigorous and fair hyperparameter tuning is still not the norm, but rather the exception in optimizer design research. Zhao et al. [2025] also examines how different optimizers perform on the pretraining tasks. However, the focus of Zhao et al. [2025] is on understanding loss structure, and the optimizers tested in the paper are shown to have slower convergence compared to AdamW, whereas the optimizers tested in this paper show small but significant improvements in convergence speed. The Algoperf competition Kasimbeg et al. [2025] evaluates different optimizers across different settings and arrives at a similar conclusion that (i) matrix-based optimizers and (ii) denoising methods such as Nesterov momentum lead to speedup over AdamW. Our works focus on the pretraining setting and investigates the effect of scaling dataset and model sizes.

### 3 Methodology

In this section, we detail the experimental design and evaluation protocol that underpin our empirical investigation. In Section 3.1, we specify the general setup for all subsequent studies. We then describe our three-phase hyperparameter-tuning framework: Phase I (in Section 3.2) performs fine-grained coordinate-descent sweeps across multiple model sizes and data-to-model ratios to identify scaling-sensitive parameters; Phase II (in Section 3.3) refines these sensitive parameters on mid-scale settings and selects the most promising optimizers; and Phase III (in Section 3.4) extrapolates hyperparameter scaling laws to the 1.2 billion-parameter regime. Together, these protocols ensure principled, fair, and reproducible comparisons across different optimizers. We present the optimal configurations found and how loss changes with respect to each hyperparameter in Appendix B and hope that this can facilitate future research. We also open-source the code (<https://github.com/marin-community/marin/tree/kaiyue/optimizers>) and the corresponding WandB runs (<https://wandb.ai/marin-community/optimizer-scaling>).

#### 3.1 General Experimental Setup

Following OLMo et al. [2025], we conduct all our experiments on a large-scale pretraining corpus composed of three publicly available datasets, tokenized with the Llama3 tokenizer: DCLM-baseline (3.8 trillion tokens, Li et al. [2025b]), StarCoder V2 Data (0.25 trillion tokens, Lozhkov et al. [2024]), and ProofPile 2 (55 billion tokens, Azerbayev et al. [2024]).

Optimizer	References	Algorithm
<b>Baseline</b>		
AdamW	Kingma and Ba [2017], Loshchilov and Hutter [2019]	Algorithm 1
<b>Variance-reduced AdamW Variants</b>		
NadamW	Dozat [2016]	Algorithm 2
Mars	Yuan et al. [2025]	Algorithm 5
Cautious	Liang et al. [2025], Wang et al. [2024]	Algorithm 7
<b>Memory-efficient Optimizers</b>		
Lion	Chen et al. [2023]	Algorithm 3
Adam-mini	Zhang et al. [2025b]	Algorithm 6
<b>Matrix-based Optimizers</b>		
Muon	Jordan et al. [2024]	Algorithm 8
Scion	Pethick et al. [2025]	Algorithm 9
Kron (PSGD)	Li [2018a, 2022]	Algorithm 10
Soap	Vyas et al. [2025]	Algorithm 11
<b>Hessian-Approximation Optimizers</b>		
Sophia	Liu et al. [2024a]	Algorithm 4

Table 1: Optimizers under study

Our benchmarks cover four model sizes derived from the architecture Touvron et al. [2023a,b], Grattafiori et al. [2024], with approximately 130M, 300M, 520M, and 1.2B parameters. Each variant uses a fixed sequence length of 4,096 and 32 transformer layers (following Liu et al. [2024b]), differing only in hidden dimension, intermediate dimension, and number of attention heads. Detailed hyperparameters are summarized in Table 2. Training is implemented in JAX and executed on TPU v5 hardware. We employ a mixed-precision scheme (parameters in fp32 and activations in bf16). For each model, we use 20 times its non-embedding parameter count to compute the Chinchilla optimal data-to-model ratio based on Hoffmann et al. [2022]. We will use  $n \times$  Chinchilla to represent training the models for  $n$  times the Chinchilla optimal number of tokens.

Model	Params	Seq Len	Hidden Dim	Inter Dim	# Layers	# Heads
Llama-130M	130M	4096	512	2048	32	8
Llama-300M	300M	4096	768	3072	32	12
Llama-520M	520M	4096	1024	4096	32	16
Llama-1.2B	1.2B	4096	1536	6144	32	24

Table 2: Detailed architecture hyperparameters for each model size we studied.

Our primary evaluation metric for the model is the language modeling loss on the English split of the C4 dataset Raffel et al. [2023], which has been shown to be a strong proxy for downstream performance Bhagia et al. [2024]. We also track downstream accuracy and bits-per-byte on the following suite of benchmarks: ARC (Easy and Challenge) Clark et al. [2018], BoolQ Clark et al. [2019], COPA Gordon et al. [2012], CommonsenseQA Talmor et al. [2018], HellaSwag Zellers et al. [2019], LAMBADA Paperno et al. [2016], OpenBookQA Mihaylov et al. [2018], PIQA Bisk et al. [2020], WSC273 Levesque et al. [2012], and Winogrande Sakaguchi et al. [2020].

Our study includes a wide range of eleven optimizers listed in Table 1. Due to the page limit, we defer the exact algorithm descriptions of these optimizers to Appendix A. We selected these eleven optimizers according to three guiding principles: (i) include widely adopted baselines such as AdamW and Lion; (ii) cover recently proposed optimizers; and (iii) when multiple methods share a similar update rule, choose a few representative algorithms. These choices ensure both breadth and depth in our comparison. We group the optimizers under study into five general classes:

1. Our baseline algorithm is AdamW, with  $m_t$  and  $v_t$  being first- and second-order momentum of gradient. AdamW’s update rule is  $w_{t+1} = w_t - \eta \frac{m_t}{\sqrt{v_t + \epsilon}} - \eta \lambda w_t$ .
2. Reducing the variance of updates is a shared motivation behind many optimizers. For example, Nesterov AdamW incorporates the Nesterov lookahead technique to estimate gradients more accurately, with the following update rule:  $w_{t+1} = w_t - \eta \frac{\beta_1 m_t + (1 - \beta_1) g_t}{\sqrt{v_t + \epsilon}} - \eta \lambda w_t$ .
3. Optimizers such as Lion aim to reduce the memory required by AdamW by only keeping the first-order momentum. Lion is observed to perform better than AdamW at a small scale Liang et al. [2025]:  $w_{t+1} = w_t - \eta \text{sign}(\beta_2 m_t + (1 - \beta_2) g_t)$ , where  $m_t$  is the first-order momentum.
4. Other optimizers like Muon leverage the matrix structure of neural networks and perform preconditioning of gradients through matrix multiplication instead of scalar multiplication. For Muon, the key operation is called Newton-Schulz:  $\text{NS}(M) = M(aM + bM^\top M + c(M \rightarrow p M)^2)$ . With appropriate  $a, b, c$ , one can prove that  $\text{NS}^{(5)}(M) \approx \arg \max_{\|O\|_{\text{op}}=1} \text{Tr}(O^\top M)$  when  $\|M\|_{\text{op}} < 1$ . Muon has the following update rule:  $w_{t+1} = w_t - \eta \text{NS}^{(5)}(\beta_2 \tilde{m}_t + (1 - \beta_2) g_t)$ . This is done for all the matrices except the token classification head and the embedding in the network.
5. Optimizers including Sophia are motivated by the famous Newton’s method and use Hessian-vector product to approximate the diagonalized version of the Hessian matrix empirically.<sup>2</sup>

### 3.2 Phase I: Fine-grained Hyperparameter Coordinate Descent

Fixed or lightly tuned baselines can severely understate an optimizer’s capability and lead to overstated speedup claims. For instance, a single tweak to the learning-rate schedule with peak learning rate 6e-4 in the GPT-3 recipe Brown et al. [2020] can produce nearly a 2× speedup (Figure 1, left). To avoid such artifacts, we perform an exhaustive, one-at-a-time sweep over each hyperparameter, identify the set of near-best configurations in each regime, and then determine which knobs truly require re-tuning as scale changes.

For each optimizer, we define a discrete grid for every hyperparameter (e.g., for AdamW, our swept hyperparameters include learning rate, weight decay, warmup steps,  $\beta_1$ ,  $\beta_2$ ,  $\epsilon$ , gradient-norm clipping, and batch size). Starting from a initial hyperparameter configuration similar to the hyperparameter configuration

<sup>2</sup>We defer the result of Sophia to Appendix B.2.

provided in the original paper proposing the optimizer, in each iteration, we hold all but one hyperparameter fixed at the current best values, then search the whole grid for that parameter, and accept the new value if the validation loss improves by more than  $\Delta_1 = 3 \times 10^{-3}$ . We repeat passes until no parameter update yields further significant gain. We perform this sweeping for 6 different settings, namely 130M, 300M, 500M at 1× Chinchilla and 130M at 2×, 4×, 8× Chinchilla. One exemplary hyperparameter optimization procedure for AdamW on a model with 300M parameters and 1× Chinchilla is shown in Table 3.

**Result of Phase I.** By the end of Phase I, we have, for each optimizer and each of the six regimes (130M, 300M, 500M at 1×; 130M at 2×, 4×, 8× Chinchilla), identified a coordinate-wise local optimum of hyperparameters—that is, the single value which, when all other knobs are held fixed, minimizes validation loss.

### 3.3 Phase II: Coordinate Descent on Scaling-Sensitive Hyperparameters

While the extensive coordinate descent guarantees coordinate-wise optimality, it is too costly to perform on larger-scale experiments. We empirically observe two crucial properties of the sweeping that allow us to simplify the descent procedure:

1. Losses are sensitive to only a subset of hyperparameters; many have little effect on performance when perturbed from their optimal values.
2. Among the sensitive hyperparameters, the optimal settings for most remain stable across scales, so tuning is needed only at smaller scales.

Building on these observations, we can simplify our hyperparameter search by identifying scaling-sensitive parameters, which (i) crucially influence the final performance, and (ii) change with respect to the model scale. We define the **approximate-optimal configuration** for a given regime as the set of all hyperparameter tuples whose final loss lies within  $\Delta_2 = 6.4\text{e-}3$  of the regime’s best-observed loss  $L_r^*$ , where  $r$  denotes a regime/setting, that is,  $r$  is a pair of choice of model size and data budget. Concretely, let  $c_h$  denotes a hyperparameter (where  $h$  is the index) and  $c$  be the tuple of all hyperparameters.

1. For each regime  $r$ , let  $\mathcal{C}_r = \{c : L(c) \leq L_r^* + \Delta_2\}$  be all the hyperparameter configurations in our coordinate descent procedure that yield approximately optimal loss.
2. A hyperparameter  $c_h$  is called **scaling-insensitive** if there exists a single value  $v_h$  such that there exists  $c \in \mathcal{C}_r$  such that  $c_h = v_h$  for every regime  $r$  in Phase I. Otherwise,  $c_h$  is **scaling-sensitive**, meaning its optimal value shifts depending on model size or data budget.

We carry forward only the scaling-sensitive hyperparameters (shown in Table 4) into Phase II, thereby focusing our next round of coordinate-descent sweeps on the hyperparameters that truly require re-tuning across scaling regimes. We then perform sweeping for another 6 different settings, namely 300M, 500M at 2×, 4×, 8× Chinchilla.

Stage	LR	WD	min lr ratio	Warmup	Max Grad Norm	Batch	Val. Loss
Init	0.008	0.1	0	1000	1	256	3.298
Round 1	0.008	0.1	0	2000	1	256	3.282
Round 2	0.008	0.1	0	2000	1	128	3.263
Best	0.008	0.1	0	2000	2	128	3.263

Table 3: Illustrative coordinate-descent steps for AdamW on the 130M 1× Chinchilla regime. Changed hyperparameter values are highlighted in red; We omitted some unchanged hyperparameters ( $\beta_1 = 0.9, \beta_2 = 0.98, \epsilon = 10^{-10}$ ).

Optimizer	Scaling-Sensitive Hyperparameters
AdamW	learning rate, warmup, weight decay, batch size
Nesterov AdamW	learning rate, warmup
Lion	learning rate, beta2
Adam-mini	learning rate, weight decay, warmup
Cautious	learning rate, batch size
Mars	learning rate, warmup, beta1
Scion	learning rate, beta1, # decay steps in WSD Hu et al. [2024]
Muon	learning rate
Soap	learning rate, warmup, block size
Kron	learning rate

Table 4: Scaling-sensitive hyperparameters identified in Phase I (Section 3.2). To reduce the memory required by Soap, we apply the parameter blocking as in Anil et al. [2021]

**Result of Phase II.** Combined with the results in Phase I, we obtain a set of near-optimal hyperparameters and their corresponding losses for 12 different settings (130M, 300M, 500M and 1×, 2×, 4×, 8× Chinchilla).

To quantify the speedup of different optimizers over the baseline AdamW, we fit how AdamW’s loss scales with data budget  $D$  for each model size  $N$  with the following functional form:  $\hat{L}_N(D) = \alpha_N D^{-B_N} + \beta_N$ . Suppose an optimizer achieves loss  $L_{\text{optimizer}}$  at data budget  $D_{\text{optimizer}}$ . We calculate the corresponding data budget needed for AdamW to achieve this loss, denoted by  $D_{\text{AdamW}}$  by finding the solution to the equation  $\hat{L}_N(D_{\text{AdamW}}) = L_{\text{optimizer}}$ . We then use  $D_{\text{AdamW}}/D_{\text{optimizer}}$  as the estimated speedup ratio.

Through this set of experiments, we observed two phenomena: (i) matrix-based optimizers consistently outperform scalar-based optimizers, but all optimizers’ speedup ratios over AdamW do not exceed 1.5×; and (ii) within matrix-based optimizers, Muon performs the best at 1-4× Chinchilla ratio but is overtaken by Soap and Kron when the Chinchilla ratio increases.

### 3.4 Phase III: Hyperparameter Scaling Law for Further Extrapolation

Having obtained optimized hyperparameter settings from Phase II (Section 3.3), we now fit a smooth scaling law that predicts the optimal value of each scaling-sensitive hyperparameter as a function of model size  $N$  and data budget  $D$ . Concretely, we model the optimal value for each scaling-sensitive hyperparameter  $h$  as:  $h(N, D) = \alpha N^{-A} D^{-B} + \beta$ , where  $A$ ,  $B$ ,  $\alpha$ , and  $\beta$  are learned coefficients.

We estimate these parameters via non-linear least-squares on the 12 observed  $(N, D, h)$  triples for each optimizer, minimizing the squared error between predicted and actual optimal hyperparameter values. To test the quality of our prediction, we ran a full Phase I sweep at  $N = 1.2\text{B}$  and Chinchilla = 1 for AdamW. Comparing the identified optimum against our fitted hyperparameters, we observe that our predicted hyperparameters yield a final loss within 3e-3 of the optimal configuration, showing that our hyperparameter scaling law can effectively predict the optimal hyperparameters. We then performed two case studies to further extrapolate our benchmarking:

1. To test the effect of scaling up model sizes, we train 1.2B models using AdamW, Nesterov AdamW, and Muon on 1 to 8× Chinchilla ratio.
2. To further test the effect of different optimizers when data-to-model ratios are high, we train 130M and 300M models using AdamW, Nesterov AdamW, Muon, and Soap on 16x Chinchilla ratio.

**Results of Phase III.** We demonstrate two potential shortcomings of Muon optimizers, which is the best optimizers in Phase I and Phase II, through experiments in this phase: (i) while Muon’s speedup persists for models up to 1.2B parameters, the speedup decreases to under 1.2×; (ii) With a 16× Chinchilla ratio, NAdamW and Soap outperform Muon on the 130M model, and Soap also surpasses Muon on the 300M model.



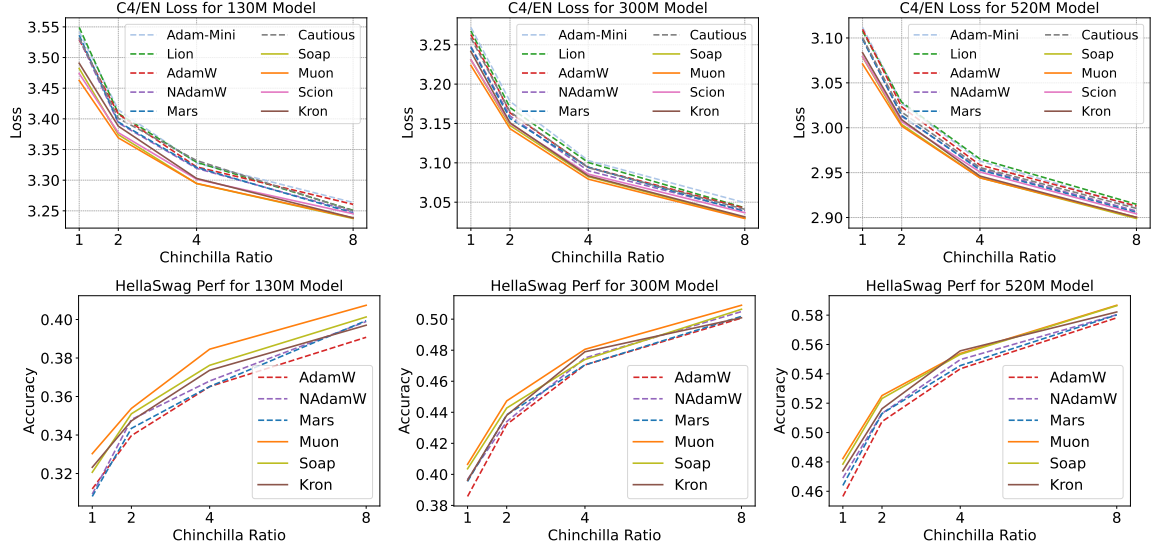


Figure 2: **Main Results For Phase I & II.** Top: We plot the validation loss on C4/EN for the experiments in Phase I and Phase II. Every point corresponds to the optimal loss achieved at the corresponding Chinchilla ratio for each optimizer. Bottom: we plot the HellaSwag performance corresponding to the selected run for a subset of optimizers: the AdamW baseline, the top 2 most performant scalar-based optimizers, and the top 3 most performant matrix-based optimizers. Analysis is deferred to Section 4.1.

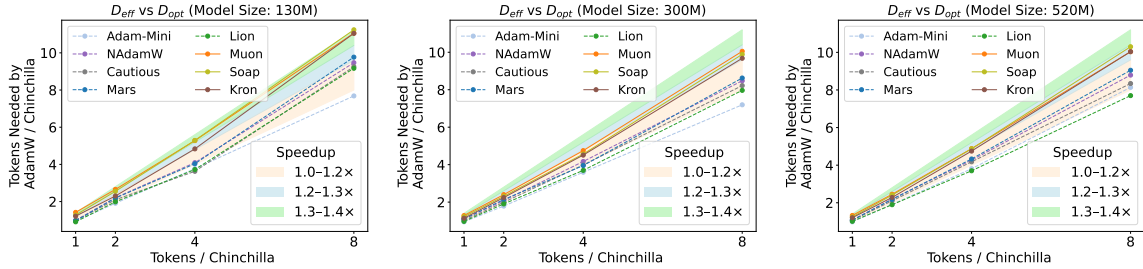


Figure 3: **Speedup of different optimizers across scale.** We estimate the speedup of different optimizers by fitting a scaling law for AdamW and then map the loss of different optimizers to the corresponding equivalent data budget. We observe that (i) The highest speedup is capped at  $1.4\times$ ; (ii) matrix-based optimizers consistently outperform scalar-based optimizers and show an increasing speedup with data budget.

## 4 Empirical Findings

### 4.1 Main Results

**Results on 0.1B–0.5B-parameter models.** Figure 2 shows the validation loss curves for the 130M, 300M, and 520M models with varying Chinchilla ratios (1 to 8) in our benchmark. We further show that HellaSwag accuracy improvements closely mirror validation-loss gains. This is consistent with prior works that show lower losses translates to better downstream accuracy Bhagia et al. [2024], Liu et al. [2025b]. Across all model scales and compute budgets, both the variance-reduced Adam variants (NAdamW, Mars, Cautious) and the matrix-based optimizers deliver speedups over the AdamW baseline. However, no method achieves  $2\times$  step-wise acceleration claimed in previous literature. We note that Soap Vyas et al. [2025] is one of the few works that conduct independent hyperparameter sweeping for the baseline, and it indeed reports a speedup closest to the actual observed improvement. Following the methodology defined in Section 3.3, we calculate the estimated speedup ratio of different optimizers in Figure 3 and the highest speedup ratio is  $1.4\times$ . From the measured speedups, three patterns stand out in this computation regime:

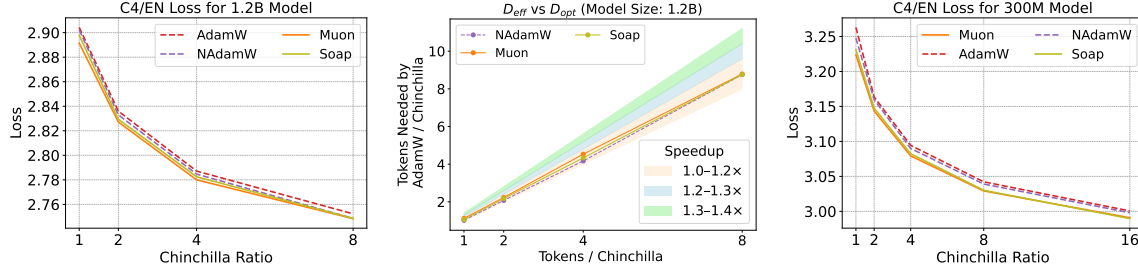


Figure 4: **Case Studies.** Left: Validation loss scaling on 1.2B model for AdamW, NAdamW, Muon and Soap. Muon and Soap still offer significant speedup over AdamW but no longer significantly speed up over NAdamW. Mid: Estimated speedup ratio with the same methodology Figure 3, we observe that Muon and Soap’s speedup decays with model size to only  $1.1\times$ . Last: Experiment with 300M  $16\times$  Chinchilla setting, Soap outperforms Muon when data-to-model ratio further increases.

1. *Matrix-based methods outperform scalar-based methods. The speedup ratio increases with data budget yet decreases with model size.* For every model size, matrix-based optimizers (Soap, Muon, Kron, Scion—solid curves) consistently drive validation loss below that of their scalar-based counterparts (dashed curves). In the base ( $1\times$  Chinchilla) compute regime, Muon performs best, but at  $8\times$  Chinchilla compute, the advantage shifts to Soap and Kron. As shown by the super-linear trend in Figure 3, the speedup ratios of these three optimizers grow with increasing data budget. To the best of authors’ knowledge, this dependency on data budget is not noted in prior works, which typically only experiment on one data-to-model ratio.

However, the greatest gains occur on the 130M model, after which the speedup decreases to roughly  $1.3\times$  for larger model sizes. We will affirm this observation further in the experiments of 1.2B models, showing that the speedup ratios of matrix-based optimizers decrease to  $1.1\times$  for the 1.2B model.

2. *Variance-reduction techniques provide a small but clear lift.* Within the scalar-based family, all variance-reduced Adam variants (NAdamW, Mars, Cautious) consistently surpass vanilla AdamW—except for a small lag at the smallest experiment. Notably, Muon combines matrix-based updates with Nesterov momentum, illustrating how variance reduction compounds with matrix adaptation can yield greater efficiency. This result is vastly different from the  $2\times$  speedup reported in some of the works and we attribute this disparity to the better-tuned baseline.
3. *Memory-efficient variants of AdamW closely track the performance of AdamW.* The two memory-efficient AdamW variants (Lion, Adam-mini)—despite their reduced auxiliary state—closely track the performance of AdamW, with a slowdown of at most 5% and sometimes even perform better than AdamW. Interestingly, Lion and Adam-mini show different scaling trends regarding model size: the disadvantage of Lion relative to AdamW widens while the disadvantage of Adam-mini over AdamW narrows.

**Results on 1.2B-parameters Models.** Using the hyperparameter scaling law we fit (Section 3.3), we scale up the model size to 1.2B to examine how the speedup of optimizers scales with model size. We observe that NAdamW, Muon, and Soap still deliver speedup over AdamW, but the speedup diminishes to  $\approx 1.1\times$  for all these optimizers (Figure 4, Left and Mid) and no longer leads to downstream improvements (Table 5). Based on this observation, we fit scaling laws for both AdamW and Muon based on the loss of 16 runs, and our scaling law predicts that Muon will result in a slightly higher loss than AdamW in the  $7B$  and  $1\times$  Chinchilla regime. We defer the fitting procedure to Appendix B.1.

**High data-to-model Ratio.** In our previous experiments, Muon is outperformed by Soap in the  $8\times$  Chinchilla regime for the 130M and 520M models. To further test this, we train three 300M models to  $16\times$  Chinchilla and verify that Muon is no longer the optimal optimizer when the data-to-model ratio increases (Figure 4, right). We conjecture that the second-order momentum maintained by Soap and Kron becomes more effective when the data-to-model ratio increases. In the long run, adaptivity to heterogeneity in parameter directions may lead to a larger speedup. We also perform similar experiments on 130M models and reach the same results (deferred to Appendix B.3).

Table 5: Benchmark performance of 1.2B models with different optimizers and Chinchilla scaling.

Optimizer	LAMBADA	OpenBook	Wino	PIQA	BoolQ	WSC273	Hella	ARC-C.	ARC-E	COPA	Avg
<b>8x Chinchilla (193B)</b>											
AdamW	67.16	41.40	64.96	76.12	68.59	82.78	67.56	43.43	74.49	85.00	67.15
NAdamW	67.84	40.20	64.80	77.15	68.10	83.52	67.40	43.34	73.61	81.00	66.70
Muon	67.53	39.80	67.09	77.09	68.81	80.95	67.67	43.34	73.53	84.00	66.98

## 4.2 Necessity of Rigorous Benchmarking

Our systematic sweeps uncover both universal optimization principles and surprising optimization-specific nuances, which call for rigorous study when designing future optimizers.

First, we find that even a superior optimizer can underperform a less advanced method when its hyperparameters are not precisely tuned. In our exhaustive grid searches, slight deviations from each optimizer’s ideal learning rate or other critical hyperparameter often lead to degradation in validation loss that is large enough to flip the ordering (Figure 5, Left). This hyperparameter sensitivity means manual selection without systematic sweeps will likely produce arbitrary ranking of the optimizers. To further demonstrate this, we plot how validation losses vary when only **one** of the hyperparameters deviates from optimal value for Muon, Soap, Mars on 520M model in  $1 \times$  Chinchilla regime (Figure 5, Mid). The ordering of the optimizers can easily flip if one chooses a sub-optimal hyperparameter.

Second, using the same hyperparameters for different optimizers do not guarantee fair comparison between optimizers. For example, while weight decay is essential across all optimizers for optimal performance, the optimal decay strength varies markedly between optimizers. As shown in (Figure 1, top right), optimal weight decay coefficients differ between the optimizers studied. We also note the optimal weight decay for Kron is larger than the conventional 0.1 and is approximately 0.5, which is crucial for Kron to outperform AdamW.

Third, early training behavior can be highly misleading. Validation-loss curves during this initial phase tend to exaggerate performance gaps (Figure 1, bottom right) and, in some cases, even reverse the eventual ranking, (Figure 5, right). Many optimizers exhibit rapid early descent followed by plateauing, meaning that assessments based solely on losses before the end of the training trajectories fail to predict final outcomes. To avoid such pitfalls, we recommend evaluating optimizers only on the final checkpoints rather than relying on intermediate checkpoints (e.g. in Liu et al. [2025a]).

We also note that the evaluation in Sophia Liu et al. [2024a] largely follows the correct procedure of the comparisons above — the peak learning rate was tuned to be optimal for the baseline. However, as the data loader in the codebase used did not fully randomize the order of the data, the optimal peak learning rate in that code base for AdamW is significantly smaller than the optimal peak learning rate for a fully randomized data loader setting. It turns out that Sophia doesn’t offer significant speedup over AdamW for models under 0.5B in our setting (Figure 7).

## 4.3 Common Phenomena Across Optimizers

Through logging the evolution of weight and gradient norms, we discovered some shared optimization phenomena across optimizers.

**Parameter norms typically track learning rate decay when there is weight decay** We observe that the parameter norms of all optimizers show a similar pattern of increase and decrease, closely aligning with the increase and decrease of the learning rate if there is a decay phase. However, the absolute values of parameter norms are very different across optimizers (Figure 6, Middle Left).

**Gradient norm increases during learning rate decay.** Across all the optimizers, the gradient norms increase during the training run. However, this increase does not lead to a loss increase. Similar to the parameter norms, the absolute values of gradient norms are not consistent across optimizers. These two phenomena are also reported in the previous work Defazio [2025], where the author also provides a theoretical explanation for both phenomena. We provide additional evidence that these two phenomena are not artifacts of the chosen optimizer AdamW but rather common phenomena across optimizers (Figure 6, Middle Right).

**Different optimizers have similar generalization behavior.** For architecture design, it has been observed that different architectures can have vastly different generalization behaviors Lu et al. [2025]. However, this is

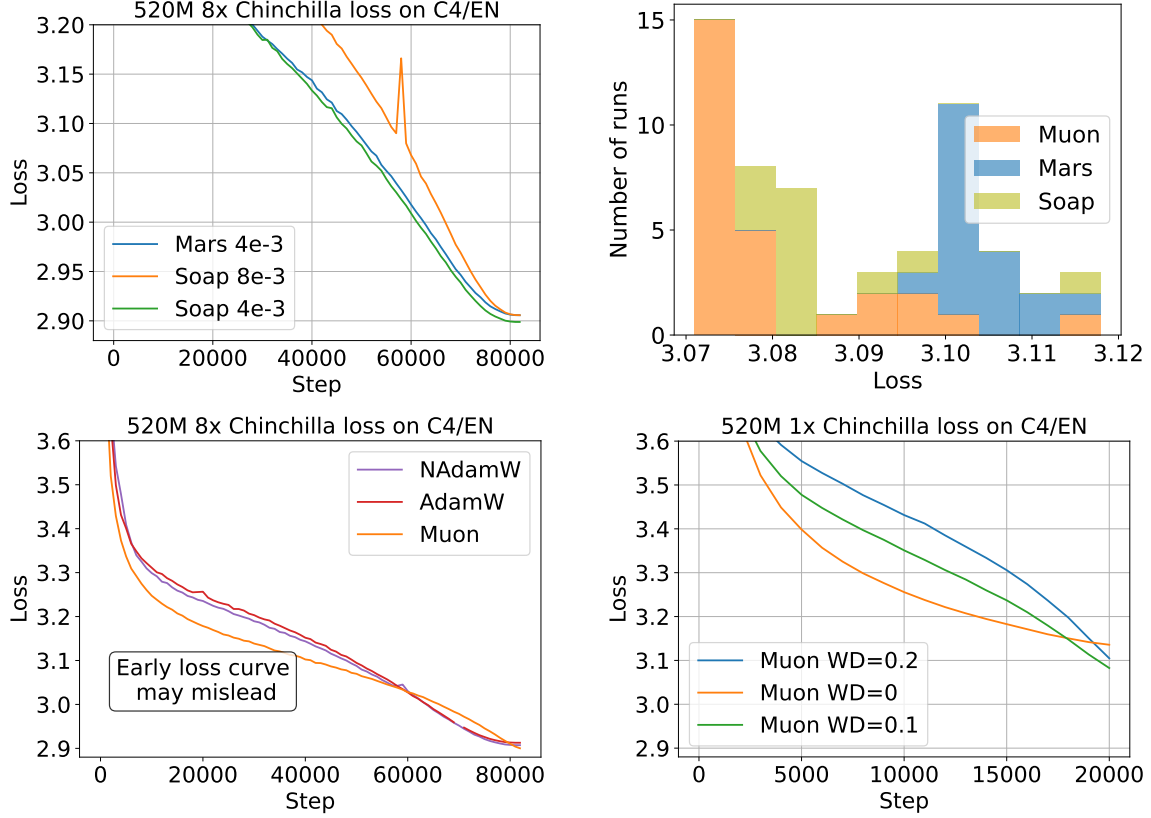


Figure 5: **Necessity of Careful Tuning.** Left:  $2\times$  the optimal learning rate diminishes Soap’s loss improvement over Mars on 520M 8x experiment; Mid: Variation of loss when only one hyperparameter differs from optimal learning rate and runs converge to within 0.02 of optimal. The order of optimizers may flip arbitrarily if rigorous tuning is missing. Right: Changing a single hyperparameter like weight decay may lead to misleading faster loss improvement but plateaus later.

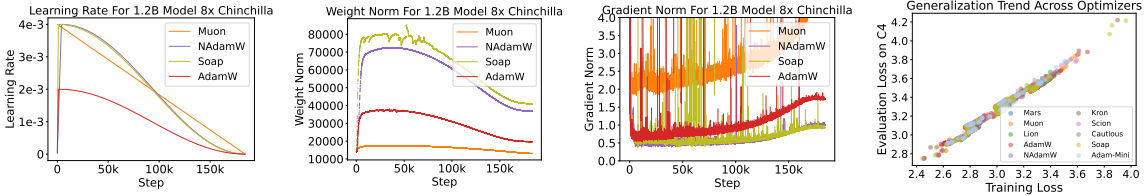


Figure 6: **Common Phenomena Across Optimizers.** Left: Learning rate used for different optimizers. Middle Left: Parameter norm of all optimizers shows a similar trend of increment and decrease, closely aligning the increasing and decaying of learning rate schedule. Middle Right: Gradient norm increases during learning rate decay. However, this increase does not lead to a loss increase. Right: The training loss and evaluation loss follows the same trend for all optimizers.

not the case for optimizers, and the evaluation losses and training losses follow roughly the same trend for all optimizers at the regimes we experimented on (Figure 6, Right).

## 5 Conclusion and Limitations

We benchmarked 11 deep learning optimizers in pretraining and found that their true gains over AdamW are much smaller than previously reported. Our results highlight three key lessons: (i) many claimed speedups

stem from insufficient hyperparameter tuning, as fair sweeps eliminate most of the apparent advantages; (ii) comparisons based on early or inconsistent evaluation can be misleading, since optimizer rankings often change over the full training trajectory; and (iii) even the best-performing alternatives provide only modest speedups, which further diminish with model scale, dropping to  $1.1\times$  at 1.2B parameters. This benchmarking study has the limitation that it does not scale to models larger than 1.2B parameters. However, we believe that evaluating optimizers on models of comparable size to prior studies is still valuable, as it reveals that insufficient tuning is a major cause of subpar speedups. Promising future work includes extending our benchmarking to larger models beyond 1.2B parameters to test whether the diminishing speedup trend persists at frontier scales. Another direction is to design optimizers whose efficiency remains stable under scaling laws, ensuring consistent benefits as model size and data budget grow.

## Acknowledgement

This work was supported by the Google TPU Research Cloud (TRC), the Stanford HAI–Google Cloud Credits Program, and NSF IIS 2211780, and is a part of the Marin Project. The authors thank Evan Walters, Omead Pooladzandi, Jiacheng You, and Zhiyuan Li for discussions and help during our research.

## References

- N. Agarwal, R. Anil, E. Hazan, T. Koren, and C. Zhang. Disentangling adaptive gradient methods from learning rates, 2020. URL <https://arxiv.org/abs/2002.11803>.
- K. Ahn, B. Xu, N. Abreu, and J. Langford. Dion: Distributed orthonormalized updates, 2025. URL <https://arxiv.org/abs/2504.05295>.
- E. AI, :, I. Shah, A. M. Polloreno, K. Stratos, P. Monk, A. Chaluvvaraju, A. Hojel, A. Ma, A. Thomas, A. Tanwer, D. J. Shah, K. Nguyen, K. Smith, M. Callahan, M. Pust, M. Parmar, P. Rushton, P. Mazarakis, R. Kapila, S. Srivastava, S. Singla, T. Romanski, Y. Vanjani, and A. Vaswani. Practical efficiency of muon for pretraining, 2025. URL <https://arxiv.org/abs/2505.02222>.
- R. Anil, V. Gupta, T. Koren, K. Regan, and Y. Singer. Scalable second order optimization for deep learning, 2021. URL <https://arxiv.org/abs/2002.09018>.
- Z. Azerbayev, H. Schoelkopf, K. Paster, M. D. Santos, S. McAleer, A. Q. Jiang, J. Deng, S. Biderman, and S. Welleck. Llemma: An open language model for mathematics, 2024. URL <https://arxiv.org/abs/2310.10631>.
- S. Becker and Y. L. Cun. Improving the convergence of back-propagation learning with second order methods. In D. S. Touretzky, G. E. Hinton, and T. J. Sejnowski, editors, *Proceedings of the 1988 Connectionist Models Summer School*, pages 29–37. San Francisco, CA: Morgan Kaufmann, 1989.
- J. Bernstein and L. Newhouse. Modular duality in deep learning, 2024a. URL <https://arxiv.org/abs/2410.21265>.
- J. Bernstein and L. Newhouse. Old optimizer, new norm: An anthology, 2024b. URL <https://arxiv.org/abs/2409.20325>.
- J. Bernstein, Y.-X. Wang, K. Azizzadenesheli, and A. Anandkumar. signSGD: Compressed optimisation for non-convex problems. In J. Dy and A. Krause, editors, *Proceedings of the 35th International Conference on Machine Learning*, volume 80 of *Proceedings of Machine Learning Research*, pages 560–569. PMLR, 10–15 Jul 2018. URL <https://proceedings.mlr.press/v80/bernstein18a.html>.
- A. Bhagia, J. Liu, A. Wettig, D. Heineman, O. Tafjord, A. H. Jha, L. Soldaini, N. A. Smith, D. Groeneveld, P. W. Koh, J. Dodge, and H. Hajishirzi. Establishing task scaling laws via compute-efficient model ladders, 2024. URL <https://arxiv.org/abs/2412.04403>.
- Y. Bisk, R. Zellers, R. Le Bras, J. Gao, and Y. Choi. Piqa: Reasoning about physical commonsense in natural language. In *Proceedings of the 34th AAAI Conference on Artificial Intelligence (AAAI)*, pages 7432–7439. AAAI Press, 2020. doi: 10.1609/aaai.v34i05.6239. URL <https://doi.org/10.1609/aaai.v34i05.6239>.
- T. B. Brown, B. Mann, N. Ryder, M. Subbiah, J. Kaplan, P. Dhariwal, A. Neelakantan, P. Shyam, G. Sastry, A. Askell, S. Agarwal, A. Herbert-Voss, G. Krueger, T. Henighan, R. Child, A. Ramesh, D. M. Ziegler, J. Wu, C. Winter, C. Hesse, M. Chen, E. Sigler, M. Litwin, S. Gray, B. Chess, J. Clark, C. Berner, S. McCandlish, A. Radford, I. Sutskever, and D. Amodei. Language models are few-shot learners, 2020. URL <https://arxiv.org/abs/2005.14165>.
- X. Chen, C. Liang, D. Huang, E. Real, K. Wang, Y. Liu, H. Pham, X. Dong, T. Luong, C.-J. Hsieh, Y. Lu, and Q. V. Le. Symbolic discovery of optimization algorithms, 2023. URL <https://arxiv.org/abs/2302.06675>.
- C. Clark, K. Lee, M.-W. Chang, T. Kwiatkowski, M. Collins, and K. Toutanova. Boolq: Exploring the surprising difficulty of natural yes/no questions. In J. Burstein, C. Doran, and T. Solorio, editors, *Proceedings of NAACL-HLT 2019*, pages 2924–2936, Minneapolis, Minnesota, June 2019. Association for Computational Linguistics. doi: 10.18653/v1/N19-1300. URL <https://aclanthology.org/N19-1300/>.

- P. Clark, I. Cowhey, O. Etzioni, T. Khot, A. Sabharwal, C. Schoenick, and O. Tafjord. Think you have solved question answering? try arc, the ai2 reasoning challenge. *CoRR*, abs/1803.05457, 2018. URL <http://arxiv.org/abs/1803.05457>.
- DeepSeek-AI, :, X. Bi, D. Chen, G. Chen, S. Chen, D. Dai, C. Deng, H. Ding, K. Dong, Q. Du, Z. Fu, H. Gao, K. Gao, W. Gao, R. Ge, K. Guan, D. Guo, J. Guo, G. Hao, Z. Hao, Y. He, W. Hu, P. Huang, E. Li, G. Li, J. Li, Y. Li, Y. K. Li, W. Liang, F. Lin, A. X. Liu, B. Liu, W. Liu, X. Liu, X. Liu, Y. Liu, H. Lu, S. Lu, F. Luo, S. Ma, X. Nie, T. Pei, Y. Piao, J. Qiu, H. Qu, T. Ren, Z. Ren, C. Ruan, Z. Sha, Z. Shao, J. Song, X. Su, J. Sun, Y. Sun, M. Tang, B. Wang, P. Wang, S. Wang, Y. Wang, Y. Wang, T. Wu, Y. Wu, X. Xie, Z. Xie, Z. Xie, Y. Xiong, H. Xu, R. X. Xu, Y. Xu, D. Yang, Y. You, S. Yu, X. Yu, B. Zhang, H. Zhang, L. Zhang, L. Zhang, M. Zhang, M. Zhang, W. Zhang, Y. Zhang, C. Zhao, Y. Zhao, S. Zhou, S. Zhou, Q. Zhu, and Y. Zou. Deepseek llm: Scaling open-source language models with longtermism, 2024. URL <https://arxiv.org/abs/2401.02954>.
- DeepSeek-AI, D. Guo, D. Yang, H. Zhang, J. Song, R. Zhang, R. Xu, Q. Zhu, S. Ma, P. Wang, X. Bi, X. Zhang, X. Yu, Y. Wu, Z. F. Wu, Z. Gou, Z. Shao, Z. Li, Z. Gao, A. Liu, B. Xue, B. Wang, B. Wu, B. Feng, C. Lu, C. Zhao, C. Deng, C. Zhang, C. Ruan, D. Dai, D. Chen, D. Ji, E. Li, F. Lin, F. Dai, F. Luo, G. Hao, G. Chen, G. Li, H. Zhang, H. Bao, H. Xu, H. Wang, H. Ding, H. Xin, H. Gao, H. Qu, H. Li, J. Guo, J. Li, J. Wang, J. Chen, J. Yuan, J. Qiu, J. Li, J. L. Cai, J. Ni, J. Liang, J. Chen, K. Dong, K. Hu, K. Gao, K. Guan, K. Huang, K. Yu, L. Wang, L. Zhang, L. Zhao, L. Wang, L. Zhang, L. Xu, L. Xia, M. Zhang, M. Zhang, M. Tang, M. Li, M. Wang, M. Li, N. Tian, P. Huang, P. Zhang, Q. Wang, Q. Chen, Q. Du, R. Ge, R. Zhang, R. Pan, R. Wang, R. J. Chen, R. L. Jin, R. Chen, S. Lu, S. Zhou, S. Chen, S. Ye, S. Wang, S. Yu, S. Zhou, S. Pan, S. S. Li, S. Zhou, S. Wu, S. Ye, T. Yun, T. Pei, T. Sun, T. Wang, W. Zeng, W. Zhao, W. Liu, W. Liang, W. Gao, W. Yu, W. Zhang, W. L. Xiao, W. An, X. Liu, X. Wang, X. Chen, X. Nie, X. Cheng, X. Liu, X. Xie, X. Liu, X. Yang, X. Li, X. Su, X. Lin, X. Q. Li, X. Jin, X. Shen, X. Chen, X. Sun, X. Wang, X. Song, X. Zhou, X. Wang, X. Shan, Y. K. Li, Y. Q. Wang, Y. X. Wei, Y. Zhang, Y. Xu, Y. Li, Y. Zhao, Y. Sun, Y. Wang, Y. Yu, Y. Zhang, Y. Shi, Y. Xiong, Y. He, Y. Piao, Y. Wang, Y. Tan, Y. Ma, Y. Liu, Y. Guo, Y. Ou, Y. Wang, Y. Gong, Y. Zou, Y. He, Y. Xiong, Y. Luo, Y. You, Y. Liu, Y. Zhou, Y. X. Zhu, Y. Xu, Y. Huang, Y. Li, Y. Zheng, Y. Zhu, Y. Ma, Y. Tang, Y. Zha, Y. Yan, Z. Z. Ren, Z. Ren, Z. Sha, Z. Fu, Z. Xu, Z. Xie, Z. Zhang, Z. Hao, Z. Ma, Z. Yan, Z. Wu, Z. Gu, Z. Zhu, Z. Liu, Z. Li, Z. Xie, Z. Song, Z. Pan, Z. Huang, Z. Xu, Z. Zhang, and Z. Zhang. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning, 2025a. URL <https://arxiv.org/abs/2501.12948>.
- DeepSeek-AI, A. Liu, B. Feng, B. Xue, B. Wang, B. Wu, C. Lu, C. Zhao, C. Deng, C. Zhang, C. Ruan, D. Dai, D. Guo, D. Yang, D. Chen, D. Ji, E. Li, F. Lin, F. Dai, F. Luo, G. Hao, G. Chen, G. Li, H. Zhang, H. Bao, H. Xu, H. Wang, H. Zhang, H. Ding, H. Xin, H. Gao, H. Li, H. Qu, J. L. Cai, J. Liang, J. Guo, J. Ni, J. Li, J. Wang, J. Chen, J. Chen, J. Yuan, J. Qiu, J. Li, J. Song, K. Dong, K. Hu, K. Gao, K. Guan, K. Huang, K. Yu, L. Wang, L. Zhang, L. Xu, L. Xia, L. Zhao, L. Wang, L. Zhang, M. Li, M. Wang, M. Zhang, M. Zhang, M. Tang, M. Li, N. Tian, P. Huang, P. Wang, P. Zhang, Q. Wang, Q. Zhu, Q. Chen, Q. Du, R. J. Chen, R. L. Jin, R. Ge, R. Zhang, R. Pan, R. Wang, R. Xu, R. Zhang, R. Chen, S. S. Li, S. Lu, S. Zhou, S. Chen, S. Wu, S. Ye, S. Ye, S. Ma, S. Wang, S. Zhou, S. Yu, S. Zhou, S. Pan, T. Wang, T. Yun, T. Pei, T. Sun, W. L. Xiao, W. Zeng, W. Zhao, W. An, W. Liu, W. Liang, W. Gao, W. Yu, W. Zhang, X. Q. Li, X. Jin, X. Wang, X. Bi, X. Liu, X. Wang, X. Shen, X. Chen, X. Zhang, X. Chen, X. Nie, X. Sun, X. Wang, X. Cheng, X. Liu, X. Xie, X. Liu, X. Yu, X. Song, X. Shan, X. Zhou, X. Yang, X. Li, X. Su, X. Lin, Y. K. Li, Y. Q. Wang, Y. X. Wei, Y. X. Zhu, Y. Zhang, Y. Xu, Y. Xu, Y. Huang, Y. Li, Y. Zhao, Y. Sun, Y. Li, Y. Wang, Y. Yu, Y. Zheng, Y. Zhang, Y. Shi, Y. Xiong, Y. He, Y. Tang, Y. Piao, Y. Wang, Y. Tan, Y. Ma, Y. Liu, Y. Guo, Y. Wu, Y. Ou, Y. Zhu, Y. Wang, Y. Gong, Y. Zou, Y. He, Y. Zha, Y. Xiong, Y. Ma, Y. Yan, Y. Luo, Y. You, Y. Liu, Y. Zhou, Z. F. Wu, Z. Z. Ren, Z. Ren, Z. Sha, Z. Fu, Z. Xu, Z. Huang, Z. Zhang, Z. Xie, Z. Zhang, Z. Hao, Z. Gou, Z. Ma, Z. Yan, Z. Shao, Z. Xu, Z. Wu, Z. Zhang, Z. Li, Z. Gu, Z. Zhu, Z. Liu, Z. Li, Z. Xie, Z. Song, Z. Gao, and Z. Pan. Deepseek-v3 technical report, 2025b. URL <https://arxiv.org/abs/2412.19437>.
- A. Defazio. Why gradients rapidly increase near the end of training, 2025. URL <https://arxiv.org/abs/2506.02285>.
- A. Defazio and K. Mishchenko. Learning-rate-free learning by d-adaptation. In A. Krause, E. Brunskill, K. Cho, B. Engelhardt, S. Sabato, and J. Scarlett, editors, *Proceedings of the 40th International Conference*

- on *Machine Learning*, volume 202 of *Proceedings of Machine Learning Research*, pages 7449–7479. PMLR, 23–29 Jul 2023. URL <https://proceedings.mlr.press/v202/defazio23a.html>.
- A. Defazio, A. Cutkosky, H. Mehta, and K. Mishchenko. Optimal linear decay learning rate schedules and further refinements, 2024a. URL <https://arxiv.org/abs/2310.07831>.
- A. Defazio, X. A. Yang, H. Mehta, K. Mishchenko, A. Khaled, and A. Cutkosky. The road less scheduled, 2024b. URL <https://arxiv.org/abs/2405.15682>.
- T. Dozat. Incorporating nesterov momentum into adam. 2016.
- J. Duchi, E. Hazan, and Y. Singer. Adaptive subgradient methods for online learning and stochastic optimization. *Journal of machine learning research*, 12(7), 2011.
- R. Eschenhagen, A. Immer, R. Turner, F. Schneider, and P. Hennig. Kronecker-factored approximate curvature for modern neural network architectures. *Advances in Neural Information Processing Systems*, 36: 33624–33655, 2023.
- R. Eschenhagen, A. Defazio, T.-H. Lee, R. E. Turner, and H.-J. M. Shi. Purifying shampoo: Investigating shampoo’s heuristics by decomposing its preconditioner, 2025. URL <https://arxiv.org/abs/2506.03595>.
- K. Everett, L. Xiao, M. Wortsman, A. A. Alemi, R. Novak, P. J. Liu, I. Gur, J. Sohl-Dickstein, L. P. Kaelbling, J. Lee, and J. Pennington. Scaling exponents across parameterizations and optimizers, 2024. URL <https://arxiv.org/abs/2407.05872>.
- P. Foret, A. Kleiner, H. Mobahi, and B. Neyshabur. Sharpness-aware minimization for efficiently improving generalization, 2021. URL <https://arxiv.org/abs/2010.01412>.
- K. Frans, S. Levine, and P. Abbeel. A stable whitening optimizer for efficient neural network training, 2025. URL <https://arxiv.org/abs/2506.07254>.
- A. Gordon, Z. Kozareva, and M. Roemmele. Semeval-2012 task 7: Choice of plausible alternatives: An evaluation of commonsense causal reasoning. In E. Agirre, J. Bos, M. Diab, S. Manandhar, Y. Marton, and D. Yuret, editors, *\*SEM 2012: The First Joint Conference on Lexical and Computational Semantics (SemEval 2012)*, pages 394–398, Montréal, Canada, June 7–8 2012. Association for Computational Linguistics. URL <https://aclanthology.org/S12-1052/>.
- A. Grattafiori, A. Dubey, A. Jauhri, A. Pandey, A. Kadian, A. Al-Dahle, A. Letman, A. Mathur, A. Schelten, A. Vaughan, A. Yang, A. Fan, A. Goyal, A. Hartshorn, A. Yang, A. Mitra, A. Srivankumar, A. Korenev, A. Hinsvark, A. Rao, A. Zhang, A. Rodriguez, A. Gregerson, A. Spataru, B. Roziere, B. Biron, B. Tang, B. Chern, C. Caucheteux, C. Nayak, C. Bi, C. Marra, C. McConnell, C. Keller, C. Touret, C. Wu, C. Wong, C. C. Ferrer, C. Nikolaidis, D. Allonsius, D. Song, D. Pintz, D. Livshits, D. Wyatt, D. Esiobu, D. Choudhary, D. Mahajan, D. Garcia-Olano, D. Perino, D. Hupkes, E. Lakomkin, E. AlBadawy, E. Lobanova, E. Dinan, E. M. Smith, F. Radenovic, F. Guzmán, F. Zhang, G. Synnaeve, G. Lee, G. L. Anderson, G. Thattai, G. Nail, G. Mialon, G. Pang, G. Cucurell, H. Nguyen, H. Korevaar, H. Xu, H. Touvron, I. Zarov, I. A. Ibarra, I. Kloumann, I. Misra, I. Evtimov, J. Zhang, J. Copet, J. Lee, J. Geffert, J. Vranes, J. Park, J. Mahadeokar, J. Shah, J. van der Linde, J. Billock, J. Hong, J. Lee, J. Fu, J. Chi, J. Huang, J. Liu, J. Wang, J. Yu, J. Bitton, J. Spisak, J. Park, J. Rocca, J. Johnstun, J. Saxe, J. Jia, K. V. Alwala, K. Prasad, K. Upasani, K. Plawiak, K. Li, K. Heafield, K. Stone, K. El-Arini, K. Iyer, K. Malik, K. Chiu, K. Bhalla, K. Lakhotia, L. Rantala-Yeary, L. van der Maaten, L. Chen, L. Tan, L. Jenkins, L. Martin, L. Madaan, L. Malo, L. Blecher, L. Landzaat, L. de Oliveira, M. Muzzi, M. Pasupuleti, M. Singh, M. Paluri, M. Kardas, M. Tsimpoukelli, M. Oldham, M. Rita, M. Pavlova, M. Kambadur, M. Lewis, M. Si, M. K. Singh, M. Hassan, N. Goyal, N. Torabi, N. Bashlykov, N. Bogoychev, N. Chatterji, N. Zhang, O. Duchenne, O. Çelebi, P. Alrassy, P. Zhang, P. Li, P. Vasic, P. Weng, P. Bhargava, P. Dubal, P. Krishnan, P. S. Koura, P. Xu, Q. He, Q. Dong, R. Srinivasan, R. Ganapathy, R. Calderer, R. S. Cabral, R. Stojnic, R. Raileanu, R. Maheswari, R. Girdhar, R. Patel, R. Sauvestre, R. Polidoro, R. Sumbaly, R. Taylor, R. Silva, R. Hou, R. Wang, S. Hosseini, S. Chennabasappa, S. Singh, S. Bell, S. S. Kim, S. Edunov, S. Nie, S. Narang, S. Raparthy, S. Shen, S. Wan, S. Bhosale, S. Zhang,



- S. Vandenhende, S. Batra, S. Whitman, S. Sootla, S. Collot, S. Gururangan, S. Borodinsky, T. Herman, T. Fowler, T. Sheasha, T. Georgiou, T. Scialom, T. Speckbacher, T. Mihaylov, T. Xiao, U. Karn, V. Goswami, V. Gupta, V. Ramanathan, V. Kerkez, V. Gonguet, V. Do, V. Vogeti, V. Albiero, V. Petrovic, W. Chu, W. Xiong, W. Fu, W. Meers, X. Martinet, X. Wang, X. Wang, X. E. Tan, X. Xia, X. Xie, X. Jia, X. Wang, Y. Goldschlag, Y. Gaur, Y. Babaei, Y. Wen, Y. Song, Y. Zhang, Y. Li, Y. Mao, Z. D. Coudert, Z. Yan, Z. Chen, Z. Papakipos, A. Singh, A. Srivastava, A. Jain, A. Kelsey, A. Shajnfeld, A. Gangidi, A. Victoria, A. Goldstand, A. Menon, A. Sharma, A. Boesenberg, A. Baeviski, A. Feinstein, A. Kallet, A. Sangani, A. Teo, A. Yunus, A. Lupu, A. Alvarado, A. Caples, A. Gu, A. Ho, A. Poulton, A. Ryan, A. Ramchandani, A. Dong, A. Franco, A. Goyal, A. Saraf, A. Chowdhury, A. Gabriel, A. Bharambe, A. Eisenman, A. Yazdan, B. James, B. Maurer, B. Leonhardi, B. Huang, B. Loyd, B. D. Paola, B. Paranjape, B. Liu, B. Wu, B. Ni, B. Hancock, B. Wasti, B. Spence, B. Stojkovic, B. Gamido, B. Montalvo, C. Parker, C. Burton, C. Mejia, C. Liu, C. Wang, C. Kim, C. Zhou, C. Hu, C.-H. Chu, C. Cai, C. Tindal, C. Feichtenhofer, C. Gao, D. Civin, D. Beaty, D. Kreymer, D. Li, D. Adkins, D. Xu, D. Testuggine, D. David, D. Parikh, D. Liskovich, D. Foss, D. Wang, D. Le, D. Holland, E. Dowling, E. Jamil, E. Montgomery, E. Presani, E. Hahn, E. Wood, E.-T. Le, E. Brinkman, E. Arcaute, E. Dunbar, E. Smothers, F. Sun, F. Kreuk, F. Tian, F. Kokkinos, F. Ozgenel, F. Caggioni, F. Kanayet, F. Seide, G. M. Florez, G. Schwarz, G. Badeer, G. Swee, G. Halpern, G. Herman, G. Sizov, Guangyi, Zhang, G. Lakshminarayanan, H. Inan, H. Shojanazeri, H. Zou, H. Wang, H. Zha, H. Habeeb, H. Rudolph, H. Suk, H. Aspegren, H. Goldman, H. Zhan, I. Damlaj, I. Molybog, I. Tufanov, I. Leontiadis, I.-E. Veliche, I. Gat, J. Weissman, J. Geboski, J. Kohli, J. Lam, J. Asher, J.-B. Gaya, J. Marcus, J. Tang, J. Chan, J. Zhen, J. Reizenstein, J. Teboul, J. Zhong, J. Jin, J. Yang, J. Cummings, J. Carvill, J. Shepard, J. McPhie, J. Torres, J. Ginsburg, J. Wang, K. Wu, K. H. U, K. Saxena, K. Khandelwal, K. Zand, K. Matosich, K. Veeraraghavan, K. Michelena, K. Li, K. Jagadeesh, K. Huang, K. Chawla, K. Huang, L. Chen, L. Garg, L. A. L. Silva, L. Bell, L. Zhang, L. Guo, L. Yu, L. Moshkovich, L. Wehrstedt, M. Khabsa, M. Avalani, M. Bhatt, M. Mankus, M. Hasson, M. Lennie, M. Reso, M. Groshev, M. Naumov, M. Lathi, M. Keneally, M. Liu, M. L. Seltzer, M. Valko, M. Restrepo, M. Patel, M. Vyatskov, M. Samvelyan, M. Clark, M. Macey, M. Wang, M. J. Hermoso, M. Metanat, M. Rastegari, M. Bansal, N. Santhanam, N. Parks, N. White, N. Bawa, N. Singhal, N. Egebo, N. Usunier, N. Mehta, N. P. Laptev, N. Dong, N. Cheng, O. Chernoguz, O. Hart, O. Salpekar, O. Kalinli, P. Kent, P. Parekh, P. Saab, P. Balaji, P. Rittner, P. Bontrager, P. Roux, P. Dollar, P. Zvyagina, P. Ratanchandani, P. Yuvraj, Q. Liang, R. Alao, R. Rodriguez, R. Ayub, R. Murthy, R. Nayani, R. Mitra, R. Parthasarathy, R. Li, R. Hogan, R. Battey, R. Wang, R. Howes, R. Rinott, S. Mehta, S. Siby, S. J. Bondu, S. Datta, S. Chugh, S. Hunt, S. Dhillon, S. Sidorov, S. Pan, S. Mahajan, S. Verma, S. Yamamoto, S. Ramaswamy, S. Lindsay, S. Lindsay, S. Feng, S. Lin, S. C. Zha, S. Patil, S. Shankar, S. Zhang, S. Zhang, S. Wang, S. Agarwal, S. Sajuyigbe, S. Chintala, S. Max, S. Chen, S. Kehoe, S. Satterfield, S. Govindaprasad, S. Gupta, S. Deng, S. Cho, S. Virk, S. Subramanian, S. Choudhury, S. Goldman, T. Remez, T. Glaser, T. Best, T. Koehler, T. Robinson, T. Li, T. Zhang, T. Matthews, T. Chou, T. Shaked, V. Vontimitta, V. Ajayi, V. Montanez, V. Mohan, V. S. Kumar, V. Mangla, V. Ionescu, V. Poenaru, V. T. Mihailescu, V. Ivanov, W. Li, W. Wang, W. Jiang, W. Bouaziz, W. Constable, X. Tang, X. Wu, X. Wang, X. Wu, X. Gao, Y. Kleinman, Y. Chen, Y. Hu, Y. Jia, Y. Qi, Y. Li, Y. Zhang, Y. Zhang, Y. Adi, Y. Nam, Yu, Wang, Y. Zhao, Y. Hao, Y. Qian, Y. Li, Y. He, Z. Rait, Z. DeVito, Z. Rosnbrick, Z. Wen, Z. Yang, Z. Zhao, and Z. Ma. The llama 3 herd of models, 2024. URL <https://arxiv.org/abs/2407.21783>.
- V. Gupta, T. Koren, and Y. Singer. Sahampoo: Preconditioned stochastic tensor optimization, 2018. URL <https://arxiv.org/abs/1802.09568>.
- J. Hoffmann, S. Borgeaud, A. Mensch, E. Buchatskaya, T. Cai, E. Rutherford, D. de Las Casas, L. A. Hendricks, J. Welbl, A. Clark, T. Hennigan, E. Noland, K. Millican, G. van den Driessche, B. Damoc, A. Guy, S. Osindero, K. Simonyan, E. Elsen, J. W. Rae, O. Vinyals, and L. Sifre. Training compute-optimal large language models, 2022. URL <https://arxiv.org/abs/2203.15556>.
- S. Hu, Y. Tu, X. Han, C. He, G. Cui, X. Long, Z. Zheng, Y. Fang, Y. Huang, W. Zhao, X. Zhang, Z. L. Thai, K. Zhang, C. Wang, Y. Yao, C. Zhao, J. Zhou, J. Cai, Z. Zhai, N. Ding, C. Jia, G. Zeng, D. Li, Z. Liu, and M. Sun. Minicpm: Unveiling the potential of small language models with scalable training strategies, 2024. URL <https://arxiv.org/abs/2404.06395>.
- M. Ivgi, O. Hinder, and Y. Carmon. Dog is sgd’s best friend: A parameter-free dynamic step size schedule, 2023. URL <https://arxiv.org/abs/2302.12022>.

- Y. Jiang, B. Neyshabur, H. Mobahi, D. Krishnan, and S. Bengio. Fantastic generalization measures and where to find them, 2019. URL <https://arxiv.org/abs/1912.02178>.
- K. Jordan, Y. Jin, V. Boza, J. You, F. Cesista, L. Newhouse, and J. Bernstein. Muon: An optimizer for hidden layers in neural networks, 2024. URL <https://kellerjordan.github.io/posts/muon/>.
- P. Kasimbeg, F. Schneider, R. Eschenhagen, J. Bae, C. S. Sastry, M. Saroufim, B. Feng, L. Wright, E. Z. Yang, Z. Nado, S. Medapati, P. Hennig, M. Rabbat, and G. E. Dahl. Accelerating neural network training: An analysis of the algoperf competition, 2025. URL <https://arxiv.org/abs/2502.15015>.
- D. P. Kingma and J. Ba. Adam: A method for stochastic optimization, 2017. URL <https://arxiv.org/abs/1412.6980>.
- T. Large, Y. Liu, M. Huh, H. Bahng, P. Isola, and J. Bernstein. Scalable optimization in the modular norm, 2024. URL <https://arxiv.org/abs/2405.14813>.
- H. J. Levesque, E. Davis, and L. Morgenstern. The winograd schema challenge. *KR*, 2012:13th, 2012.
- H. Li, W. Zheng, J. Hu, Q. Wang, H. Zhang, Z. Wang, S. Xuyang, Y. Fan, S. Zhou, X. Zhang, and D. Jiang. Predictable scale: Part i – optimal hyperparameter scaling law in large language model pretraining, 2025a. URL <https://arxiv.org/abs/2503.04715>.
- J. Li, A. Fang, G. Smyrnis, M. Ivgi, M. Jordan, S. Gadre, H. Bansal, E. Guha, S. Keh, K. Arora, S. Garg, R. Xin, N. Muennighoff, R. Heckel, J. Mercat, M. Chen, S. Gururangan, M. Wortsman, A. Albalak, Y. Bitton, M. Nezhurina, A. Abbas, C.-Y. Hsieh, D. Ghosh, J. Gardner, M. Kilian, H. Zhang, R. Shao, S. Pratt, S. Sanyal, G. Ilharco, G. Daras, K. Marathe, A. Gokaslan, J. Zhang, K. Chandu, T. Nguyen, I. Vasiljevic, S. Kakade, S. Song, S. Sanghavi, F. Faghri, S. Oh, L. Zettlemoyer, K. Lo, A. El-Nouby, H. Pouransari, A. Toshev, S. Wang, D. Groeneveld, L. Soldaini, P. W. Koh, J. Jitsev, T. Kollar, A. G. Dimakis, Y. Carmon, A. Dave, L. Schmidt, and V. Shankar. Datacomp-lm: In search of the next generation of training sets for language models, 2025b. URL <https://arxiv.org/abs/2406.11794>.
- X. Li. Black box lie group preconditioners for sgd, 2022. URL <https://arxiv.org/abs/2211.04422>.
- X.-L. Li. Preconditioned stochastic gradient descent. *IEEE Transactions on Neural Networks and Learning Systems*, 29(5):1454–1466, May 2018a. ISSN 2162-2388. doi: 10.1109/tnnls.2017.2672978. URL <http://dx.doi.org/10.1109/TNNLS.2017.2672978>.
- X.-L. Li. Preconditioner on matrix lie group for sgd, 2018b. URL <https://arxiv.org/abs/1809.10232>.
- K. Liang, L. Chen, B. Liu, and Q. Liu. Cautious optimizers: Improving training with one line of code, 2025. URL <https://arxiv.org/abs/2411.16085>.
- H. Liu, Z. Li, D. Hall, P. Liang, and T. Ma. Sophia: A scalable stochastic second-order optimizer for language model pre-training, 2024a. URL <https://arxiv.org/abs/2305.14342>.
- J. Liu, J. Su, X. Yao, Z. Jiang, G. Lai, Y. Du, Y. Qin, W. Xu, E. Lu, J. Yan, Y. Chen, H. Zheng, Y. Liu, S. Liu, B. Yin, W. He, H. Zhu, Y. Wang, J. Wang, M. Dong, Z. Zhang, Y. Kang, H. Zhang, X. Xu, Y. Zhang, Y. Wu, X. Zhou, and Z. Yang. Muon is scalable for llm training, 2025a. URL <https://arxiv.org/abs/2502.16982>.
- L. Liu, H. Jiang, P. He, W. Chen, X. Liu, J. Gao, and J. Han. On the variance of the adaptive learning rate and beyond, 2021. URL <https://arxiv.org/abs/1908.03265>.
- Q. Liu, X. Zheng, N. Muennighoff, G. Zeng, L. Dou, T. Pang, J. Jiang, and M. Lin. Regmix: Data mixture as regression for language model pre-training, 2025b. URL <https://arxiv.org/abs/2407.01492>.
- Y. Liu, Z. Liu, and J. Gore. Focus: First order concentrated updating scheme, 2025c. URL <https://arxiv.org/abs/2501.12243>.
- Z. Liu, C. Zhao, F. Iandola, C. Lai, Y. Tian, I. Fedorov, Y. Xiong, E. Chang, Y. Shi, R. Krishnamoorthi, L. Lai, and V. Chandra. Mobilellm: Optimizing sub-billion parameter language models for on-device use cases, 2024b. URL <https://arxiv.org/abs/2402.14905>.

- Z. Liu, Y. Liu, E. J. Michaud, J. Gore, and M. Tegmark. Physics of skill learning, 2025d. URL <https://arxiv.org/abs/2501.12391>.
- I. Loshchilov and F. Hutter. Decoupled weight decay regularization, 2019. URL <https://arxiv.org/abs/1711.05101>.
- A. Lozhkov, R. Li, L. B. Allal, F. Cassano, J. Lamy-Poirier, N. Tazi, A. Tang, D. Pykhtar, J. Liu, Y. Wei, T. Liu, M. Tian, D. Kocetkov, A. Zucker, Y. Belkada, Z. Wang, Q. Liu, D. Abulkhanov, I. Paul, Z. Li, W.-D. Li, M. Risdal, J. Li, J. Zhu, T. Y. Zhuo, E. Zheltonozhskii, N. O. O. Dade, W. Yu, L. Krauß, N. Jain, Y. Su, X. He, M. Dey, E. Abati, Y. Chai, N. Muennighoff, X. Tang, M. Oblokulov, C. Akiki, M. Marone, C. Mou, M. Mishra, A. Gu, B. Hui, T. Dao, A. Zebaze, O. Dehaene, N. Patry, C. Xu, J. McAuley, H. Hu, T. Scholak, S. Paquet, J. Robinson, C. J. Anderson, N. Chapados, M. Patwary, N. Tajbakhsh, Y. Jernite, C. M. Ferrandis, L. Zhang, S. Hughes, T. Wolf, A. Guha, L. von Werra, and H. de Vries. Starcoder 2 and the stack v2: The next generation, 2024. URL <https://arxiv.org/abs/2402.19173>.
- X. Lu, Y. Zhao, S. Wei, S. Wang, B. Qin, and T. Liu. How does sequence modeling architecture influence base capabilities of pre-trained language models? exploring key architecture design principles to avoid base capabilities degradation, 2025. URL <https://arxiv.org/abs/2505.18522>.
- L. Luo, Y. Xiong, Y. Liu, and X. Sun. Adaptive gradient methods with dynamic bound of learning rate, 2019. URL <https://arxiv.org/abs/1902.09843>.
- Y. Luo, X. Ren, Z. Zheng, Z. Jiang, X. Jiang, and Y. You. Came: Confidence-guided adaptive memory efficient optimization, 2023. URL <https://arxiv.org/abs/2307.02047>.
- C. Ma, W. Gong, M. Scetbon, and E. Meeds. Swan: Sgd with normalization and whitening enables stateless llm training, 2025. URL <https://arxiv.org/abs/2412.13148>.
- J. Martens and R. Grosse. Optimizing neural networks with kronecker-factored approximate curvature, 2020. URL <https://arxiv.org/abs/1503.05671>.
- T. Mihaylov, P. Clark, T. Khot, and A. Sabharwal. Can a suit of armor conduct electricity? a new dataset for open book question answering. In E. Riloff, D. Chiang, J. Hockenmaier, and J. Tsujii, editors, *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 2381–2391, Brussels, Belgium, Oct–Nov 2018. Association for Computational Linguistics. doi: 10.18653/v1/D18-1260. URL <https://aclanthology.org/D18-1260/>.
- K. Mishchenko and A. Defazio. Prodigy: An expeditiously adaptive parameter-free learner, 2024. URL <https://arxiv.org/abs/2306.06101>.
- I.-V. Modoranu, M. Safaryan, G. Malinovsky, E. Kurtić, T. Robert, P. Richtárik, and D. Alistarh. Microadam: Accurate adaptive optimization with low space overhead and provable convergence. *Advances in Neural Information Processing Systems*, 37:1–43, 2024.
- D. Morwani, I. Shapira, N. Vyas, E. Malach, S. Kakade, and L. Janson. A new perspective on shampoo’s preconditioner, 2024. URL <https://arxiv.org/abs/2406.17748>.
- Y. Nesterov. A method for solving the convex programming problem with convergence rate  $O(1/k^2)$ . In *Dokl akad nauk Sssr*, volume 269, page 543, 1983.
- T. OLMO, P. Walsh, L. Soldaini, D. Groeneveld, K. Lo, S. Arora, A. Bhagia, Y. Gu, S. Huang, M. Jordan, N. Lambert, D. Schwenk, O. Tafjord, T. Anderson, D. Atkinson, F. Brahman, C. Clark, P. Dasigi, N. Dziri, M. Guerquin, H. Ivison, P. W. Koh, J. Liu, S. Malik, W. Merrill, L. J. V. Miranda, J. Morrison, T. Murray, C. Nam, V. Pyatkin, A. Rangapur, M. Schmitz, S. Skjongsberg, D. Wadden, C. Wilhelm, M. Wilson, L. Zettlemoyer, A. Farhadi, N. A. Smith, and H. Hajishirzi. 2 olmo 2 furious, 2025. URL <https://arxiv.org/abs/2501.00656>.
- M. Pagliardini, P. Ablin, and D. Grangier. The ademamix optimizer: Better, faster, older, 2024. URL <https://arxiv.org/abs/2409.03137>.

- D. Paperno, G. Kruszewski, A. Lazaridou, N. Q. Pham, R. Bernardi, S. Pezzelle, M. Baroni, G. Boleda, and R. Fernández. The LAMBADA dataset: Word prediction requiring a broad discourse context. In K. Erk and N. A. Smith, editors, *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1525–1534, Berlin, Germany, August 2016. Association for Computational Linguistics. doi: 10.18653/v1/P16-1144. URL <https://aclanthology.org/P16-1144/>.
- T. Pethick, W. Xie, K. Antonakopoulos, Z. Zhu, A. Silveti-Falls, and V. Cevher. Training deep learning models with norm-constrained Imos, 2025. URL <https://arxiv.org/abs/2502.07529>.
- C. Raffel, N. Shazeer, A. Roberts, K. Lee, S. Narang, M. Matena, Y. Zhou, W. Li, and P. J. Liu. Exploring the limits of transfer learning with a unified text-to-text transformer, 2023. URL <https://arxiv.org/abs/1910.10683>.
- S. J. Reddi, S. Kale, and S. Kumar. On the convergence of adam and beyond. In *International Conference on Learning Representations*, 2018. URL <https://openreview.net/forum?id=ryQu7f-RZ>.
- H. Robbins and S. Monro. A stochastic approximation method. *The annals of mathematical statistics*, pages 400–407, 1951.
- K. Sakaguchi, R. Le Bras, C. Bhagavatula, and Y. Choi. Winogrande: An adversarial winograd schema challenge at scale. In *Proceedings of the 34th AAAI Conference on Artificial Intelligence (AAAI)*, pages 8732–8740. AAAI Press, 2020. doi: 10.1609/aaai.v34i05.6399. URL <https://doi.org/10.1609/aaai.v34i05.6399>.
- T. Schaul, S. Zhang, and Y. LeCun. No more pesky learning rates, 2013. URL <https://arxiv.org/abs/1206.1106>.
- R. M. Schmidt, F. Schneider, and P. Hennig. Descending through a crowded valley - benchmarking deep learning optimizers, 2021. URL <https://arxiv.org/abs/2007.01547>.
- N. Shazeer and M. Stern. Adafactor: Adaptive learning rates with sublinear memory cost, 2018. URL <https://arxiv.org/abs/1804.04235>.
- I. Sutskever, J. Martens, G. Dahl, and G. Hinton. On the importance of initialization and momentum in deep learning. In S. Dasgupta and D. McAllester, editors, *Proceedings of the 30th International Conference on Machine Learning*, volume 28 of *Proceedings of Machine Learning Research*, pages 1139–1147, Atlanta, Georgia, USA, 17–19 Jun 2013. PMLR. URL <https://proceedings.mlr.press/v28/sutskever13.html>.
- A. Talmor, J. Herzig, N. Lourie, and J. Berant. Commonsenseqa: A question answering challenge targeting commonsense knowledge. *CoRR*, abs/1811.00937, 2018. URL <http://arxiv.org/abs/1811.00937>.
- S. Taniguchi, K. Harada, G. Minegishi, Y. Oshima, S. C. Jeong, G. Nagahara, T. Iiyama, M. Suzuki, Y. Iwasawa, and Y. Matsuo. Adopt: Modified adam can converge with any  $\beta_2$  with the optimal rate, 2024. URL <https://arxiv.org/abs/2411.02853>.
- K. Team, Y. Bai, Y. Bao, G. Chen, J. Chen, N. Chen, R. Chen, Y. Chen, Y. Chen, Y. Chen, Z. Chen, J. Cui, H. Ding, M. Dong, A. Du, C. Du, D. Du, Y. Du, Y. Fan, Y. Feng, K. Fu, B. Gao, H. Gao, P. Gao, T. Gao, X. Gu, L. Guan, H. Guo, J. Guo, H. Hu, X. Hao, T. He, W. He, W. He, C. Hong, Y. Hu, Z. Hu, W. Huang, Z. Huang, Z. Huang, T. Jiang, Z. Jiang, X. Jin, Y. Kang, G. Lai, C. Li, F. Li, H. Li, M. Li, W. Li, Y. Li, Y. Li, Z. Li, Z. Li, H. Lin, X. Lin, Z. Lin, C. Liu, C. Liu, H. Liu, J. Liu, J. Liu, L. Liu, S. Liu, T. Y. Liu, T. Liu, W. Liu, Y. Liu, Y. Liu, Y. Liu, Y. Liu, Z. Liu, E. Lu, L. Lu, S. Ma, X. Ma, Y. Ma, S. Mao, J. Mei, X. Men, Y. Miao, S. Pan, Y. Peng, R. Qin, B. Qu, Z. Shang, L. Shi, S. Shi, F. Song, J. Su, Z. Su, X. Sun, F. Sung, H. Tang, J. Tao, Q. Teng, C. Wang, D. Wang, F. Wang, H. Wang, J. Wang, J. Wang, J. Wang, S. Wang, S. Wang, Y. Wang, Y. Wang, Y. Wang, Y. Wang, Y. Wang, Z. Wang, Z. Wang, Z. Wang, C. Wei, Q. Wei, W. Wu, X. Wu, Y. Wu, C. Xiao, X. Xie, W. Xiong, B. Xu, J. Xu, J. Xu, L. H. Xu, L. Xu, S. Xu, W. Xu, X. Xu, Y. Xu, Z. Xu, J. Yan, Y. Yan, X. Yang, Y. Yang, Z. Yang, Z. Yang, Z. Yang, H. Yao, X. Yao, W. Ye, Z. Ye, B. Yin, L. Yu, E. Yuan, H. Yuan, M. Yuan, H. Zhan, D. Zhang, H. Zhang, W. Zhang, X. Zhang, Y. Zhang,

- Y. Zhang, Y. Zhang, Y. Zhang, Y. Zhang, Y. Zhang, Z. Zhang, H. Zhao, Y. Zhao, H. Zheng, S. Zheng, J. Zhou, X. Zhou, Z. Zhou, Z. Zhu, W. Zhuang, and X. Zu. Kimi k2: Open agentic intelligence, 2025. URL <https://arxiv.org/abs/2507.20534>.
- T. Tieleman. Lecture 6.5-rmsprop: Divide the gradient by a running average of its recent magnitude. *COURSERA: Neural networks for machine learning*, 4(2):26, 2012.
- H. Touvron, T. Lavril, G. Izacard, X. Martinet, M.-A. Lachaux, T. Lacroix, B. Rozière, N. Goyal, E. Hambro, F. Azhar, A. Rodriguez, A. Joulin, E. Grave, and G. Lample. Llama: Open and efficient foundation language models, 2023a. URL <https://arxiv.org/abs/2302.13971>.
- H. Touvron, L. Martin, K. Stone, P. Albert, A. Almahairi, Y. Babaei, N. Bashlykov, S. Batra, P. Bhargava, S. Bhosale, D. Bikel, L. Blecher, C. C. Ferrer, M. Chen, G. Cucurull, D. Esiobu, J. Fernandes, J. Fu, W. Fu, B. Fuller, C. Gao, V. Goswami, N. Goyal, A. Hartshorn, S. Hosseini, R. Hou, H. Inan, M. Kardas, V. Kerkez, M. Khabsa, I. Kloumann, A. Korenev, P. S. Koura, M.-A. Lachaux, T. Lavril, J. Lee, D. Liskovich, Y. Lu, Y. Mao, X. Martinet, T. Mihaylov, P. Mishra, I. Molybog, Y. Nie, A. Poulton, J. Reizenstein, R. Rungta, K. Saladi, A. Schelten, R. Silva, E. M. Smith, R. Subramanian, X. E. Tan, B. Tang, R. Taylor, A. Williams, J. X. Kuan, P. Xu, Z. Yan, I. Zarov, Y. Zhang, A. Fan, M. Kambadur, S. Narang, A. Rodriguez, R. Stojnic, S. Edunov, and T. Scialom. Llama 2: Open foundation and fine-tuned chat models, 2023b. URL <https://arxiv.org/abs/2307.09288>.
- N. Vyas, D. Morwani, R. Zhao, M. Kwun, I. Shapira, D. Brandfonbrener, L. Janson, and S. Kakade. Soap: Improving and stabilizing shampoo using adam, 2025. URL <https://arxiv.org/abs/2409.11321>.
- J. Wang, M. Wang, Z. Zhou, J. Yan, W. E, and L. Wu. The sharpness disparity principle in transformers for accelerating language model pre-training, 2025. URL <https://arxiv.org/abs/2502.19002>.
- S. Wang, A. Liu, J. Xiao, H. Liu, Y. Yang, C. Xu, Q. Pu, S. Zheng, W. Zhang, and J. Li. Cadam: Confidence-based optimization for online learning, 2024. URL <https://arxiv.org/abs/2411.19647>.
- K. Wen, Z. Li, J. Wang, D. Hall, P. Liang, and T. Ma. Understanding warmup-stable-decay learning rates: A river valley loss landscape perspective, 2024. URL <https://arxiv.org/abs/2410.05192>.
- X. Xie, P. Zhou, H. Li, Z. Lin, and S. Yan. Adan: Adaptive nesterov momentum algorithm for faster optimizing deep models, 2024. URL <https://arxiv.org/abs/2208.06677>.
- Z. Xie, X. Wang, H. Zhang, I. Sato, and M. Sugiyama. Adaptive inertia: Disentangling the effects of adaptive learning rate and momentum, 2022. URL <https://arxiv.org/abs/2006.15815>.
- A. Yang, A. Li, B. Yang, B. Zhang, B. Hui, B. Zheng, B. Yu, C. Gao, C. Huang, C. Lv, C. Zheng, D. Liu, F. Zhou, F. Huang, F. Hu, H. Ge, H. Wei, H. Lin, J. Tang, J. Yang, J. Tu, J. Zhang, J. Yang, J. Yang, J. Zhou, J. Zhou, J. Lin, K. Dang, K. Bao, K. Yang, L. Yu, L. Deng, M. Li, M. Xue, M. Li, P. Zhang, P. Wang, Q. Zhu, R. Men, R. Gao, S. Liu, S. Luo, T. Li, T. Tang, W. Yin, X. Ren, X. Wang, X. Zhang, X. Ren, Y. Fan, Y. Su, Y. Zhang, Y. Zhang, Y. Wan, Y. Liu, Z. Wang, Z. Cui, Z. Zhang, Z. Zhou, and Z. Qiu. Qwen3 technical report, 2025. URL <https://arxiv.org/abs/2505.09388>.
- G. Yang. Scaling limits of wide neural networks with weight sharing: Gaussian process behavior, gradient independence, and neural tangent kernel derivation, 2020a. URL <https://arxiv.org/abs/1902.04760>.
- G. Yang. Tensor programs iii: Neural matrix laws. *arXiv preprint arXiv:2009.10685*, 2020b.
- G. Yang. Tensor programs ii: Neural tangent kernel for any architecture. *arXiv preprint arXiv:2006.14548*, 2020c.
- G. Yang. Tensor programs i: Wide feedforward or recurrent neural networks of any architecture are gaussian processes, 2021. URL <https://arxiv.org/abs/1910.12478>.
- G. Yang and E. Littwin. Tensor programs ivb: Adaptive optimization in the infinite-width limit, 2023. URL <https://arxiv.org/abs/2308.01814>.

- G. Yang, E. J. Hu, I. Babuschkin, S. Sidor, X. Liu, D. Farhi, N. Ryder, J. Pachocki, W. Chen, and J. Gao. Tensor programs v: Tuning large neural networks via zero-shot hyperparameter transfer, 2022. URL <https://arxiv.org/abs/2203.03466>.
- G. Yang, J. B. Simon, and J. Bernstein. A spectral condition for feature learning, 2024. URL <https://arxiv.org/abs/2310.17813>.
- Z. Yao, A. Gholami, S. Shen, M. Mustafa, K. Keutzer, and M. Mahoney. Adahessian: An adaptive second order optimizer for machine learning. *Proceedings of the AAAI Conference on Artificial Intelligence*, 35(12): 10665–10673, May 2021. doi: 10.1609/aaai.v35i12.17275. URL <https://ojs.aaai.org/index.php/AAAI/article/view/17275>.
- Y. You, I. Gitman, and B. Ginsburg. Large batch training of convolutional networks, 2017. URL <https://arxiv.org/abs/1708.03888>.
- Y. You, J. Li, S. Reddi, J. Hseu, S. Kumar, S. Bhojanapalli, X. Song, J. Demmel, K. Keutzer, and C.-J. Hsieh. Large batch optimization for deep learning: Training bert in 76 minutes, 2020. URL <https://arxiv.org/abs/1904.00962>.
- H. Yuan, Y. Liu, S. Wu, X. Zhou, and Q. Gu. Mars: Unleashing the power of variance reduction for training large models, 2025. URL <https://arxiv.org/abs/2411.10438>.
- M. Zaheer, S. Reddi, D. Sachan, S. Kale, and S. Kumar. Adaptive methods for nonconvex optimization. *Advances in neural information processing systems*, 31, 2018.
- M. D. Zeiler. Adadelata: An adaptive learning rate method, 2012. URL <https://arxiv.org/abs/1212.5701>.
- R. Zellers, A. Holtzman, Y. Bisk, A. Farhadi, and Y. Choi. Hellaswag: Can a machine really finish your sentence? In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics (ACL)*, pages 4791–4800, 2019. doi: 10.18653/V1/P19-1472. URL <https://doi.org/10.18653/V1/P19-1472>.
- H. Zhang, D. Morwani, N. Vyas, J. Wu, D. Zou, U. Ghai, D. Foster, and S. Kakade. How does critical batch size scale in pre-training?, 2025a. URL <https://arxiv.org/abs/2410.21676>.
- M. R. Zhang, J. Lucas, G. Hinton, and J. Ba. Lookahead optimizer: k steps forward, 1 step back, 2019. URL <https://arxiv.org/abs/1907.08610>.
- Y. Zhang, C. Chen, Z. Li, T. Ding, C. Wu, D. P. Kingma, Y. Ye, Z.-Q. Luo, and R. Sun. Adam-mini: Use fewer learning rates to gain more, 2025b. URL <https://arxiv.org/abs/2406.16793>.
- J. Zhao, Z. Zhang, B. Chen, Z. Wang, A. Anandkumar, and Y. Tian. Galore: Memory-efficient llm training by gradient low-rank projection, 2024. URL <https://arxiv.org/abs/2403.03507>.
- R. Zhao, D. Morwani, D. Brandfonbrener, N. Vyas, and S. Kakade. Deconstructing what makes a good optimizer for language models, 2025. URL <https://arxiv.org/abs/2407.07972>.
- H. Zhu, Z. Zhang, W. Cong, X. Liu, S. Park, V. Chandra, B. Long, D. Z. Pan, Z. Wang, and J. Lee. Apollo: Sgd-like memory, adamw-level performance, 2025. URL <https://arxiv.org/abs/2412.05270>.
- J. Zhuang, T. Tang, Y. Ding, S. Tatikonda, N. Dvornek, X. Papademetris, and J. S. Duncan. Adabelief optimizer: Adapting stepsizes by the belief in observed gradients, 2020. URL <https://arxiv.org/abs/2010.07468>.

## A Optimizer Definitions

In this section, we present the algorithm for each optimizer we evaluated.

Throughout this section, we use the following notation:  $w_t$  for model parameters,  $g_t$  for gradients at step  $t$ ,  $\eta$  for learning rate,  $\lambda$  for weight decay,  $\beta_1, \beta_2$  for moment decay rates,  $\epsilon$  for numerical stability,  $g_{\text{norm}}$  for gradient norm, and  $m, v$  for first and second moments. All operations are element-wise unless specified.

---

### Algorithm 1 AdamW

---

**Hyperparameters:**  $\beta_1, \beta_2, \epsilon, \eta, \lambda, g_{\text{norm}}$

**State:**  $m, v$

**Update Rule:**

$$\begin{aligned}\hat{g}_t &= g_t \max\{1, \frac{g_{\text{norm}}}{\|g_t\|_2}\} \\ m_t &= \beta_1 m_{t-1} + (1 - \beta_1) \hat{g}_t, \\ v_t &= \beta_2 v_{t-1} + (1 - \beta_2) \hat{g}_t^2, \\ \hat{m}_t &= \frac{m_t}{1 - \beta_1^t}, \quad \hat{v}_t = \frac{v_t}{1 - \beta_2^t}, \\ w_{t+1} &= w_t - \eta \frac{\hat{m}_t}{\sqrt{\hat{v}_t} + \epsilon} - \eta \lambda w_t.\end{aligned}$$


---

---

### Algorithm 2 Nesterov AdamW

---

**Hyperparameters:**  $\beta_1, \beta_2, \epsilon, \eta, \lambda, g_{\text{norm}}$

**State:**  $m, v$

**Update Rule:**

$$\begin{aligned}\hat{g}_t &= g_t \max\{1, \frac{g_{\text{norm}}}{\|g_t\|_2}\} \\ m_t &= \beta_1 m_{t-1} + (1 - \beta_1) \hat{g}_t, \\ v_t &= \beta_2 v_{t-1} + (1 - \beta_2) \hat{g}_t^2, \\ \tilde{m}_t &= \beta_1 m_t + (1 - \beta_1) \hat{g}_t, \\ \hat{m}_t &= \frac{\tilde{m}_t}{1 - \beta_1^t}, \quad \hat{v}_t = \frac{v_t}{1 - \beta_2^t}, \\ w_{t+1} &= w_t - \eta \frac{\hat{m}_t}{\sqrt{\hat{v}_t} + \epsilon} - \eta \lambda w_t.\end{aligned}$$


---

---

**Algorithm 3** Lion

---

**Hyperparameters:**  $\beta_1, \beta_2, \eta, \epsilon, \lambda, g_{\text{norm}}$

**State:**  $m$

**Update Rule:**

$$\begin{aligned}\hat{g}_t &= g_t \max\{1, \frac{g_{\text{norm}}}{\|g_t\|_2}\} \\ \hat{m}_t &= \beta_1 m_{t-1} + (1 - \beta_1) \hat{g}_t, \\ m_{t+1} &= \beta_2 m_{t-1} + (1 - \beta_2) \hat{g}_t, \\ w_{t+1} &= w_t - \eta \text{sign}(\hat{m}_t) - \eta \lambda w_t.\end{aligned}$$

---

**Algorithm 4** Sophia-H

---

**Hyperparameters:**  $\{\eta_t\}_{t=1}^T, \lambda, k, \beta_1, \beta_2, \gamma, \epsilon$

**State:**  $m_0 = 0, h_{1-k} = 0, \theta_1$

**For**  $t = 1, \dots, T$ :

$$\begin{aligned}g_t &= \nabla_{\theta} L_t(\theta_t) \\ m_t &= \beta_1 m_{t-1} + (1 - \beta_1) g_t \\ \text{If } t \bmod k &= 1: \\ r &\sim \{\pm 1\}^d, \quad v = g_t \cdot r, \quad u = \nabla_{\theta} v, \\ \hat{h} &= r \odot u \\ h_t &= \beta_2 h_{t-k} + (1 - \beta_2) \hat{h} \\ \text{Else: } h_t &= h_{t-1} \\ \theta_t &= \theta_t - \eta_t \lambda \theta_t \\ \theta_{t+1} &= \theta_t - \eta_t \text{clip}\left(\frac{m_t}{\max(\gamma h_t, \epsilon)}, 1\right)\end{aligned}$$



---

**Algorithm 5** MARS

---

**Hyperparameters:**  $\beta_1, \beta_2, \gamma, \epsilon, \eta, \lambda, g_{\text{norm}}$

**State:**  $m, v, g_{t-1}$

**Update Rule:**

$$\begin{aligned}c_t &= g_t + \gamma \frac{\beta_1}{1 - \beta_1} (g_t - g_{t-1}), \\ \hat{c}_t &= c_t \max\{1, \frac{g_{\text{norm}}}{\|c_t\|_2}\}, \\ m_t &= \beta_1 m_{t-1} + (1 - \beta_1) \hat{c}_t, \\ v_t &= \beta_2 v_{t-1} + (1 - \beta_2) \hat{c}_t^2, \\ \hat{m}_t &= \frac{m_t}{1 - \beta_1^t}, \quad \hat{v}_t = \frac{v_t}{1 - \beta_2^t}, \\ w_{t+1} &= w_t - \eta \frac{\hat{m}_t}{\sqrt{\hat{v}_t} + \epsilon} - \eta \lambda w_t.\end{aligned}$$

---

**Algorithm 6** Adam-mini

---

**Hyperparameters:**  $\beta_1, \beta_2, \epsilon, \eta, \lambda, g_{\text{norm}}$

**State:**  $m$  (with the same shape as  $w$ ),  $v$  (one scalar for each block)

**Setup:**

1. Partition all parameters into `param_blocks`: please refer to Zhang et al. [2025b] for the exact partition scheme.
2. We will use  $w_{t,b}$  and  $m_{t,b}$  to denote the parameters in block  $b$  at step  $t$

**Update**

$$\begin{aligned}\hat{g}_t &= g_t \max\{1, \frac{g_{\text{norm}}}{\|g_t\|_2}\} \\ m_t &= \beta_1 m_{t-1} + (1 - \beta_1) \hat{g}_t, \\ \hat{m}_t &= \frac{m_t}{1 - \beta_1^t},\end{aligned}$$

**for each block**  $b \in \text{param\_blocks}$

$$\begin{aligned}v_{t,b} &= \beta_2 v_{t-1,b} + (1 - \beta_2) \text{mean}(g_{t,b}^2) \\ \hat{v}_{t,b} &= \frac{v_{t,b}}{1 - \beta_2^t}.\end{aligned}$$

$$w_{t,b} = w_{t-1,b} - \eta \hat{m}_{t,b} / (\sqrt{\hat{v}_{t,b}} + \epsilon) - \eta \lambda w_{t-1,b}.$$

---

**Algorithm 7** Cautious

---

**Hyperparameters:**  $\beta_1, \beta_2, \epsilon, \eta, \lambda, g_{\text{norm}}$

**State:**  $m, v$

**Update Rule:**

$$\begin{aligned}\hat{g}_t &= g_t \max\{1, \frac{g_{\text{norm}}}{\|g_t\|_2}\} \\ m_t &= \beta_1 m_{t-1} + (1 - \beta_1) \hat{g}_t, \\ v_t &= \beta_2 v_{t-1} + (1 - \beta_2) \hat{g}_t^2, \\ \hat{m}_t &= \frac{m_t}{1 - \beta_1^t}, \quad \hat{v}_t = \frac{v_t}{1 - \beta_2^t}, \\ u_t &= \frac{\hat{m}_t}{\sqrt{\hat{v}_t + \epsilon}}, \quad s_t = \mathbb{I}(u_t \cdot \hat{g}_t > 0), \\ \hat{u}_t &= \frac{u_t \cdot s_t}{\text{mean}(s_t)}, \quad w_{t+1} = w_t - \eta \hat{u}_t - \eta \lambda w_t.\end{aligned}$$

---

**Algorithm 8** Muon

---

**Hyperparameters:**  $\beta, \eta, \epsilon, \beta_1, \beta_2, \epsilon_{\text{Adam}}, \eta_{\text{Adam}}, \lambda, g_{\text{norm}}$

**State:**  $m$

**Update Rule For Weights in LM Head, Embedding, and LayerNorm:** Same as AdamW

**Update Rule For Matrices in Transformer Layer :**

$$\begin{aligned}\hat{g}_t &= g_t \max\{1, \frac{g_{\text{norm}}}{\|g_t\|_2}\} \\ m_t &= \beta m_{t-1} + \hat{g}_t, \\ u &= \beta m_t + \hat{g}_t, \\ u &= \text{NewtonSchulz}(u, \text{steps} = 5), \\ s &= \sqrt{\max(1, \frac{\text{rows}(u)}{\text{cols}(u)})}, \\ u &= s u, \\ w_{t+1} &= w_t - \eta u - \eta \lambda w_t.\end{aligned}$$

**Newton–Schulz Orthogonalization (Operating on Matrices):**

$$\begin{aligned}X &\leftarrow \frac{X}{\|X\| + \epsilon}, \quad \text{transpose} \leftarrow (\text{rows}(X) > \text{cols}(X)), \\ \text{if transpose: } X &\leftarrow X^\top, \\ \text{for } i = 1 \dots 5 : A &= X X^\top, B = 3.4445 A - 4.7750 A^2 + 2.0315 A^3, \\ X &\leftarrow 3.4445 X + B X, \\ \text{if transpose: } X &\leftarrow X^\top, \quad \text{return } X.\end{aligned}$$

---

**Algorithm 9** Scion

---

**Hyperparameters:**  $\beta, \eta, \epsilon, \beta_1, \beta_2, \eta_{\text{SignGD}}, g_{\text{norm}}$

**State:**  $m$

**Update Rule For Matrices in LM Head and Embedding:**

$$\begin{aligned}\hat{g}_t &= g_t \max\left\{1, \frac{g_{\text{norm}}}{\|g_t\|_2}\right\} \\ m_t &= \beta_1 m_{t-1} + (1 - \beta_1) \hat{g}_t, \\ w_{t+1} &= w_t - \eta_{\text{SignGD}} \text{sign}(\hat{m}_t).\end{aligned}$$

**Update Rule For Matrices in Transformer Layer :** Same as Muon

---

---

**Algorithm 10** PSGD Kron

**Hyperparameters:**  $\beta_1, \eta, \lambda, \epsilon, g_{\text{norm}}, \text{normalize\_grads}, \text{partition\_grads\_into\_blocks}, \text{merge\_small\_dims}, \text{block\_size}, \text{target\_merged\_dim\_size}, p_{\text{upd}}(t)$ ,

**Setup:**

If `merge_small_dims` is True, then try merging small dimensions into a single dimension with size `target_merged_dim_size` greedily.

If `partition_grads_into_blocks` is True, then partition all parameters into `block_size`  $\times$  `block_size` blocks.

Define  $\text{unfold}_i$  as the function that unfolds all dimensions except the  $i$ -th dimension into a single dimension and  $\text{fold}_i$  as the inverse function.

**State (per block  $\ell$ ):** Denote  $w^{(\ell)}$  as the parameters in block  $\ell$  and assume it has shape  $d_1 \times d_2 \times \dots \times d_n$ .

$\mu^{(\ell)}$ , with the same shape as  $w^{(\ell)}$ ,

$Q_i^{(\ell)}$  ( $i = 1, \dots, n$ ), a lower-triangular matrix with shape  $d_i \times d_i$

if  $g_{\text{norm}} > 0$ :  $\hat{g}_t = g_t \max\{1, \frac{g_{\text{norm}}}{\|g_t\|_2}\}$ , else  $\hat{g}_t = g_t$

if `normalize_grads`:  $\hat{g}_t = \frac{\hat{g}_t}{\|\hat{g}_t\|_2 + \epsilon}$ , else  $\hat{g}_t = \hat{g}_t \mu_t \leftarrow \beta_1 \mu_{t-1} + (1 - \beta_1) \hat{g}_t$ ,  $\hat{\mu}_t \leftarrow \frac{\mu_t}{1 - \beta_1^t}$ .

**Balance check:**  $\text{bal\_ctr} \leftarrow \text{bal\_ctr} + 1$ . If  $\text{bal\_ctr} \geq 100$ , then  $Q_i^{(\ell)} \leftarrow \text{balance}(Q_i^{(\ell)})$ ,  $\text{bal\_ctr} \leftarrow 0$ .

**Balance function:**  $\text{balance}(Q_i^{(\ell)})$  performs the following steps:

1. Compute the maximum absolute value for each row/column of  $Q_i^{(\ell)}$ :  $\text{norms}_i = \max_{j,k} |Q_i^{(\ell)}[j, k]|$
2. Calculate the geometric mean:  $\text{gmean}_i = (\prod_{j=1}^n \text{norms}_i[j])^{\frac{1}{d_i}}$
3. Scale each element:  $Q_i^{(\ell)} \leftarrow Q_i^{(\ell)} \cdot \frac{\text{gmean}_i}{\text{norms}_i}$
4. Return the balanced  $Q_i^{(\ell)}$

**Preconditioner update:** with probability  $p_{\text{upd}}(t)$ ,

1. **Random probe:**  $V^{(\ell)} \leftarrow \text{tree\_random\_like}(g_t^{(\ell)})$ .
2. **Dampen:**  $\hat{\mu}_t^{(\ell)} \leftarrow \hat{\mu}_t^{(\ell)} + \epsilon \cdot \text{mean}(|\hat{\mu}_t^{(\ell)}|) \cdot V^{(\ell)}$
3. **Conjugate sketch ( $B$ ):**

$$X^{(0)} = V^{(\ell)}, \quad X^{(i)} = \text{fold}_i \left( \left( Q_i^{(\ell)} \right)^{-T} \text{unfold}_i(X^{(i-1)}) \right), \quad B = X^{(n)}.$$

4. **Pre-sketch ( $A$ ):**

$$Y^{(0)} = \hat{\mu}^{(\ell)}, \quad Y^{(i)} = \text{fold}_i \left( Q_i^{(\ell)} \text{unfold}_i(Y^{(i-1)}) \right), \quad A = Y^{(n)}.$$

5. For each  $i$  in  $1, \dots, n$ ,

$$M_i = \text{unfold}_i(A), \quad C_i = \text{unfold}_i(B),$$

$$T_1 = M_i M_i^T, \quad T_2 = C_i C_i^T, \quad s = \max_{u,v} |(T_1 + T_2)_{u,v}|,$$

$$Q_i^{(\ell)} \leftarrow Q_i^{(\ell)} - \alpha \frac{T_1 - T_2}{s} Q_i^{(\ell)}.$$

$$G^{(0)} = \hat{\mu}^{(\ell)}, \quad G^{(i)} = \text{fold}_i \left( Q_i^{(\ell)} \text{unfold}_i(G^{(i-1)}) \right), \quad \tilde{g}_t^{(\ell)} = G^{(d)}.$$

**(Weight-decay & update)**

28

$$\tilde{g}_t^{(\ell)} \leftarrow \tilde{g}_t^{(\ell)} + \lambda w_{t-1}^{(\ell)}, \quad w_t^{(\ell)} \leftarrow w_{t-1}^{(\ell)} - \eta \tilde{g}_t^{(\ell)}.$$

---

**Algorithm 11** SOAP

---

**Hyperparameters:**  $\beta_1, \beta_2, \mu, k, \epsilon, \text{block\_size}, g_{\text{norm}}$

**State:**  $m, v, G_A, G_B, Q_A, Q_B$

Partition all parameters into  $\text{block\_size} \times \text{block\_size}$  block

**Update Rule For Each Block:**

$$\begin{aligned}\hat{g}_t &= g_t \max\{1, \frac{g_{\text{norm}}}{\|g_t\|_2}\} \\ \hat{g}_t &= Q_A \hat{g}_t Q_B, \\ m_t &= \beta_1 m_{t-1} + (1 - \beta_1) \hat{g}_t, \quad v_t = \beta_2 v_{t-1} + (1 - \beta_2) \hat{g}_t^2, \\ \hat{m}_t &= \frac{m_t}{1 - \beta_1^t}, \quad \hat{v}_t = \frac{v_t}{1 - \beta_2^t}, \\ w_{t+1} &= w_t - \eta_t Q_A^\top \left( \frac{\hat{m}_t}{\sqrt{\hat{v}_t} + \epsilon} \right) Q_B^\top, \\ G_A &= \mu G_A + (1 - \mu) \hat{g}_t \hat{g}_t^\top + \epsilon I, \quad G_B = \mu G_B + (1 - \mu) \hat{g}_t^\top \hat{g}_t + \epsilon I, \\ \text{if } t \bmod k = 0 : \quad & Q_A = \text{QR}(G_A Q_A), \quad Q_B = \text{QR}(G_B Q_B).\end{aligned}$$

## B Omitted Experiments

### B.1 Scaling Law

Based on Appendices C and D, we fitted our scaling law for the 1.2B run of Muon, NAdamW and AdamW, we round the fitted value to our grid of hyperparameter and used hyperparameters are shown in Appendix E. After training the models, we fitted a scaling law for Muon and AdamW of the following form:

$$L(N, D) = \alpha N^{-A} + \beta D^{-B} + \gamma \quad (1)$$

The fitted values are

- For AdamW,  $\alpha = 21.4289, A = 0.1555, \beta = 276.4235, B = 0.2804, \gamma = 1.7324$ , with a RMS error of  $3 \times 10^{-3}$ .
- For Muon  $\alpha = 32.7458, A = 0.1864, \beta = 59.0221, B = 0.2074, \gamma = 1.8063$ , with a RMS error of  $5 \times 10^{-3}$ .

This scaling law predicts that when the parameter scale reaches 7B, Muon will actually result in a higher loss compared to AdamW in  $1 \times$  Chinchilla regime.

### B.2 Sophia Experiments

We performed a Phase I experiment with Sophia, with detailed hyperparameter settings shown in Appendix C. We found that Sophia tends to underperform AdamW in smaller compute regimes and eventually slightly outperforms AdamW when either the model size or the data size increases.

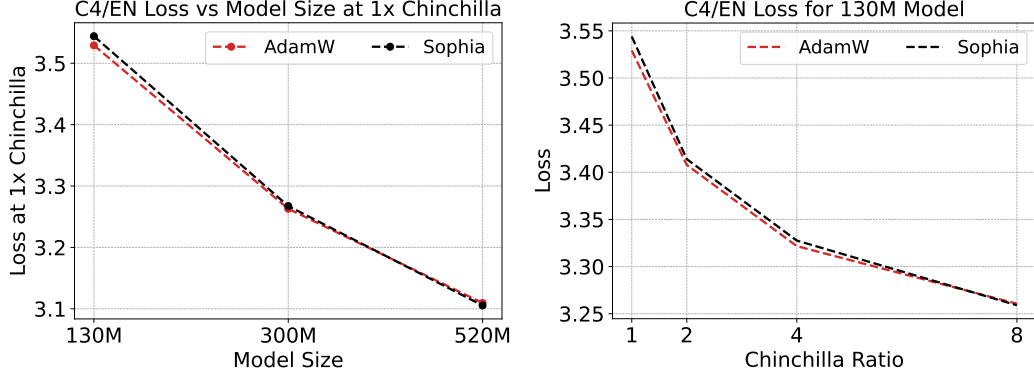


Figure 7: **Sophia Experiments.** Left: Loss curve of Sophia and AdamW in  $1\times$  Chinchilla setting. Right: Loss curve of Sophia and AdamW for 130M model size.

### B.3 High Data-to-model Ratio

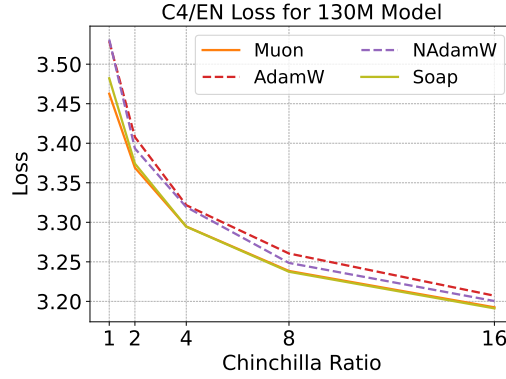


Figure 8: **More Case Studies.** Experiment with 130M  $16\times$  Chinchilla setting, SOAP outperforms Muon in the overtraining setting.

### B.4 Evaluation Performance

Table 6: Evaluation Performance for Mars, Model Size = 130m

DATA SIZE	2B	5B	10B	21B
PERFORMANCE METRIC				
LAMBADA OPENAI	0.327	0.377	0.402	0.451
OPENBOOKQA	0.274	0.304	0.298	0.302
WINOGRANDE	0.521	0.512	0.526	0.522
PIQA	0.625	0.652	0.662	0.678
BOOLQ	0.616	0.504	0.569	0.557
WSC273	0.557	0.560	0.579	0.590
HELLASWAG 0SHOT	0.308	0.343	0.365	0.399
ARC CHALLENGE	0.225	0.258	0.272	0.247
ARC EASY	0.462	0.517	0.532	0.544
COPA	0.650	0.690	0.660	0.710
FINAL C4 LOSS	3.537	3.396	3.323	3.247

Table 7: Evaluation Performance for Mars, Model Size = 300m

DATA SIZE	6B	12B	24B	48B
PERFORMANCE METRIC				
LAMBADA OPENAI	0.450	0.493	0.514	0.538
OPENBOOKQA	0.300	0.312	0.328	0.328
WINOGRANDE	0.534	0.557	0.549	0.561
PIQA	0.680	0.686	0.695	0.720
BOOLQ	0.480	0.533	0.612	0.608
WSC273	0.586	0.608	0.667	0.667
HELLASWAG 0SHOT	0.395	0.437	0.471	0.502
ARC CHALLENGE	0.253	0.269	0.298	0.307
ARC EASY	0.547	0.580	0.608	0.616
COPA	0.700	0.740	0.720	0.730
FINAL C4 LOSS	3.249	3.158	3.097	3.040

Table 8: Evaluation Performance for Mars, Model Size = 520m

DATA SIZE	10B	21B	42B	85B
PERFORMANCE METRIC				
LAMBADA OPENAI	0.513	0.548	0.576	0.601
OPENBOOKQA	0.330	0.336	0.366	0.378
WINOGRANDE	0.548	0.559	0.586	0.601
PIQA	0.704	0.717	0.725	0.735
BOOLQ	0.559	0.598	0.625	0.631
WSC273	0.601	0.696	0.736	0.780
HELLASWAG 0SHOT	0.462	0.513	0.546	0.579
ARC CHALLENGE	0.294	0.317	0.333	0.351
ARC EASY	0.611	0.648	0.665	0.676
COPA	0.730	0.760	0.750	0.750
FINAL C4 LOSS	3.101	3.015	2.955	2.906

Table 9: Evaluation Performance for Muon, Model Size = 130m

DATA SIZE	2B	5B	10B	21B
PERFORMANCE METRIC				
LAMBADA OPENAI	0.358	0.396	0.431	0.456
OPENBOOKQA	0.286	0.288	0.318	0.300
WINOGRANDE	0.515	0.518	0.530	0.541
PIQA	0.640	0.660	0.670	0.678
BOOLQ	0.571	0.504	0.583	0.452
WSC273	0.553	0.564	0.579	0.582
HELLASWAG 0SHOT	0.330	0.354	0.381	0.407
ARC CHALLENGE	0.231	0.244	0.269	0.276
ARC EASY	0.493	0.511	0.543	0.559
COPA	0.650	0.720	0.730	0.660
FINAL C4 LOSS	3.464	3.369	3.296	3.240

Table 10: Evaluation Performance for Muon, Model Size = 300m

DATA SIZE	6B	12B	24B	48B
PERFORMANCE METRIC				
LAMBADA OPENAI	0.459	0.492	0.529	0.550
OPENBOOKQA	0.302	0.310	0.330	0.346
WINOGRANDE	0.504	0.534	0.548	0.570
PIQA	0.679	0.689	0.695	0.714
BOOLQ	0.482	0.555	0.599	0.570
WSC273	0.615	0.641	0.659	0.656
HELLASWAG 0SHOT	0.406	0.447	0.481	0.509
ARC CHALLENGE	0.270	0.289	0.290	0.316
ARC EASY	0.564	0.591	0.620	0.638
COPA	0.670	0.730	0.700	0.700
FINAL C4 LOSS	3.224	3.143	3.079	3.032

Table 11: Evaluation Performance for Muon, Model Size = 520m

DATA SIZE	10B	21B	42B	85B
PERFORMANCE METRIC				
LAMBADA OPENAI	0.520	0.559	0.597	0.611
OPENBOOKQA	0.328	0.340	0.342	0.368
WINOGRANDE	0.560	0.567	0.594	0.590
PIQA	0.713	0.718	0.730	0.736
BOOLQ	0.590	0.555	0.633	0.613
WSC273	0.659	0.667	0.762	0.736
HELLASWAG 0SHOT	0.482	0.525	0.554	0.587
ARC CHALLENGE	0.300	0.318	0.351	0.348
ARC EASY	0.623	0.640	0.671	0.674
COPA	0.730	0.760	0.730	0.780
FINAL C4 LOSS	3.073	3.002	2.945	2.900

Table 12: Evaluation Performance for Lion, Model Size = 130m

DATA SIZE	2B	5B	10B	21B
PERFORMANCE METRIC				
LAMBADA OPENAI	0.316	0.380	0.417	0.450
OPENBOOKQA	0.280	0.296	0.312	0.292
WINOGRANDE	0.503	0.518	0.519	0.521
PIQA	0.629	0.649	0.663	0.680
BOOLQ	0.603	0.507	0.419	0.489
WSC273	0.564	0.538	0.590	0.623
HELLASWAG 0SHOT	0.304	0.343	0.366	0.402
ARC CHALLENGE	0.229	0.272	0.262	0.272
ARC EASY	0.464	0.500	0.536	0.539
COPA	0.700	0.670	0.730	0.730
FINAL C4 LOSS	3.552	3.409	3.331	3.252



Table 13: Evaluation Performance for Lion, Model Size = 300m

DATA SIZE	6B	12B	24B	48B
PERFORMANCE METRIC				
LAMBADA OPENAI	0.428	0.471	0.509	0.535
OPENBOOKQA	0.306	0.320	0.326	0.346
WINOGRANDE	0.506	0.530	0.532	0.564
PIQA	0.676	0.686	0.699	0.713
BOOLQ	0.486	0.490	0.609	0.593
WSC273	0.579	0.615	0.608	0.667
HELLASWAG 0SHOT	0.386	0.428	0.466	0.502
ARC CHALLENGE	0.267	0.296	0.294	0.322
ARC EASY	0.535	0.586	0.600	0.618
COPA	0.690	0.730	0.710	0.760
FINAL C4 LOSS	3.268	3.170	3.100	3.046

Table 14: Evaluation Performance for Lion, Model Size = 520m

DATA SIZE	10B	21B	42B	85B
PERFORMANCE METRIC				
LAMBADA OPENAI	0.507	0.543	0.569	0.597
OPENBOOKQA	0.336	0.334	0.356	0.358
WINOGRANDE	0.560	0.548	0.581	0.590
PIQA	0.699	0.713	0.726	0.740
BOOLQ	0.600	0.546	0.612	0.631
WSC273	0.608	0.667	0.685	0.700
HELLASWAG 0SHOT	0.463	0.509	0.545	0.571
ARC CHALLENGE	0.303	0.320	0.340	0.371
ARC EASY	0.603	0.638	0.664	0.673
COPA	0.710	0.750	0.740	0.770
FINAL C4 LOSS	3.108	3.029	2.965	2.915

Table 15: Evaluation Performance for NAdamW, Model Size = 130m

DATA SIZE	2B	5B	10B	21B
PERFORMANCE METRIC				
LAMBADA OPENAI	0.325	0.383	0.413	0.446
OPENBOOKQA	0.284	0.286	0.308	0.322
WINOGRANDE	0.515	0.507	0.546	0.507
PIQA	0.633	0.652	0.667	0.682
BOOLQ	0.591	0.569	0.459	0.557
WSC273	0.535	0.560	0.612	0.597
HELLASWAG 0SHOT	0.312	0.348	0.368	0.399
ARC CHALLENGE	0.239	0.245	0.242	0.270
ARC EASY	0.465	0.503	0.518	0.547
COPA	0.680	0.700	0.650	0.720
FINAL C4 LOSS	3.531	3.394	3.319	3.251

Table 16: Evaluation Performance for NAdamW, Model Size = 300m

DATA SIZE	6B	12B	24B	48B
PERFORMANCE METRIC				
LAMBADA OPENAI	0.437	0.477	0.514	0.543
OPENBOOKQA	0.310	0.332	0.316	0.348
WINOGRANDE	0.520	0.536	0.549	0.564
PIQA	0.681	0.693	0.708	0.712
BOOLQ	0.539	0.564	0.520	0.612
WSC273	0.582	0.637	0.645	0.700
HELLASWAG 0SHOT	0.397	0.434	0.475	0.505
ARC CHALLENGE	0.264	0.275	0.292	0.307
ARC EASY	0.553	0.589	0.605	0.632
COPA	0.710	0.730	0.700	0.750
FINAL C4 LOSS	3.248	3.160	3.090	3.039

Table 17: Evaluation Performance for NAdamW, Model Size = 520m

DATA SIZE	10B	21B	42B	85B
PERFORMANCE METRIC				
LAMBADA OPENAI	0.500	0.551	0.577	0.599
OPENBOOKQA	0.332	0.350	0.358	0.342
WINOGRANDE	0.553	0.571	0.594	0.597
PIQA	0.695	0.719	0.731	0.738
BOOLQ	0.613	0.606	0.639	0.618
WSC273	0.637	0.689	0.718	0.714
HELLASWAG 0SHOT	0.469	0.513	0.550	0.580
ARC CHALLENGE	0.294	0.322	0.342	0.354
ARC EASY	0.605	0.646	0.657	0.674
COPA	0.700	0.720	0.760	0.780
FINAL C4 LOSS	3.100	3.013	2.955	2.907

Table 18: Evaluation Performance for Kron, Model Size = 130m

DATA SIZE	2B	5B	10B	21B
PERFORMANCE METRIC				
LAMBADA OPENAI	0.341	0.387	0.426	0.451
OPENBOOKQA	0.268	0.278	0.296	0.314
WINOGRANDE	0.520	0.522	0.515	0.543
PIQA	0.636	0.662	0.662	0.683
BOOLQ	0.533	0.530	0.575	0.557
WSC273	0.564	0.564	0.553	0.604
HELLASWAG 0SHOT	0.323	0.347	0.374	0.397
ARC CHALLENGE	0.235	0.259	0.256	0.281
ARC EASY	0.487	0.519	0.543	0.572
COPA	0.700	0.690	0.710	0.710
FINAL C4 LOSS	3.492	3.389	3.307	3.239

Table 19: Evaluation Performance for Kron, Model Size = 300m

DATA SIZE	6B	12B	24B	48B
PERFORMANCE METRIC				
LAMBADA OPENAI	0.429	0.485	0.516	0.548
OPENBOOKQA	0.310	0.314	0.314	0.326
WINOGRANDE	0.527	0.534	0.542	0.571
PIQA	0.675	0.690	0.706	0.709
BOOLQ	0.487	0.475	0.610	0.611
WSC273	0.546	0.619	0.659	0.718
HELLASWAG 0SHOT	0.396	0.438	0.479	0.501
ARC CHALLENGE	0.262	0.285	0.298	0.323
ARC EASY	0.551	0.586	0.618	0.632
COPA	0.670	0.720	0.720	0.750
FINAL C4 LOSS	3.244	3.151	3.083	3.031

Table 20: Evaluation Performance for Kron, Model Size = 520m

DATA SIZE	10B	21B	42B	85B
PERFORMANCE METRIC				
LAMBADA OPENAI	0.512	0.549	0.569	0.603
OPENBOOKQA	0.338	0.332	0.350	0.352
WINOGRANDE	0.536	0.567	0.575	0.601
PIQA	0.721	0.713	0.724	0.739
BOOLQ	0.529	0.560	0.537	0.559
WSC273	0.648	0.663	0.696	0.762
HELLASWAG 0SHOT	0.474	0.516	0.556	0.582
ARC CHALLENGE	0.287	0.312	0.344	0.353
ARC EASY	0.616	0.648	0.662	0.681
COPA	0.710	0.720	0.790	0.750
FINAL C4 LOSS	3.084	3.009	2.946	2.900

Table 21: Evaluation Performance for Scion, Model Size = 130m

DATA SIZE	2B	5B	10B	21B
PERFORMANCE METRIC				
LAMBADA OPENAI	0.355	0.382	0.427	0.446
OPENBOOKQA	0.282	0.300	0.288	0.314
WINOGRANDE	0.515	0.502	0.502	0.535
PIQA	0.641	0.655	0.675	0.684
BOOLQ	0.574	0.518	0.532	0.474
WSC273	0.505	0.531	0.553	0.608
HELLASWAG 0SHOT	0.323	0.354	0.374	0.401
ARC CHALLENGE	0.247	0.241	0.276	0.264
ARC EASY	0.501	0.513	0.541	0.557
COPA	0.650	0.660	0.660	0.680
FINAL C4 LOSS	3.477	3.379	3.302	3.246

Table 22: Evaluation Performance for Scion, Model Size = 300m

DATA SIZE	6B	12B	24B	48B
PERFORMANCE METRIC				
LAMBADA OPENAI	0.445	0.489	0.511	0.540
OPENBOOKQA	0.306	0.308	0.326	0.358
WINOGRANDE	0.531	0.550	0.555	0.566
PIQA	0.684	0.689	0.701	0.712
BOOLQ	0.576	0.610	0.574	0.585
WSC273	0.597	0.612	0.685	0.696
HELLASWAG 0SHOT	0.406	0.442	0.478	0.506
ARC CHALLENGE	0.272	0.285	0.298	0.316
ARC EASY	0.564	0.592	0.625	0.633
COPA	0.680	0.730	0.700	0.740
FINAL C4 LOSS	3.232	3.152	3.086	3.039

Table 23: Evaluation Performance for Scion, Model Size = 520m

DATA SIZE	10B	21B	42B	85B
PERFORMANCE METRIC				
LAMBADA OPENAI	0.521	0.552	0.575	0.601
OPENBOOKQA	0.338	0.362	0.366	0.374
WINOGRANDE	0.569	0.551	0.605	0.586
PIQA	0.704	0.717	0.736	0.743
BOOLQ	0.575	0.615	0.620	0.641
WSC273	0.674	0.692	0.736	0.736
HELLASWAG 0SHOT	0.481	0.520	0.549	0.581
ARC CHALLENGE	0.289	0.314	0.346	0.354
ARC EASY	0.625	0.647	0.667	0.686
COPA	0.730	0.730	0.750	0.770
FINAL C4 LOSS	3.080	3.007	2.952	2.904

Table 24: Evaluation Performance for Cautious, Model Size = 130m

DATA SIZE	2B	5B	10B	21B
PERFORMANCE METRIC				
LAMBADA OPENAI	0.314	0.385	0.420	0.458
OPENBOOKQA	0.298	0.280	0.308	0.318
WINOGRANDE	0.515	0.520	0.519	0.508
PIQA	0.638	0.652	0.664	0.669
BOOLQ	0.580	0.512	0.512	0.550
WSC273	0.560	0.527	0.593	0.604
HELLASWAG 0SHOT	0.314	0.347	0.371	0.401
ARC CHALLENGE	0.239	0.244	0.261	0.270
ARC EASY	0.476	0.508	0.532	0.544
COPA	0.690	0.720	0.710	0.690
FINAL C4 LOSS	3.535	3.403	3.334	3.253

Table 25: Evaluation Performance for Cautious, Model Size = 300m

DATA SIZE	6B	12B	24B	48B
PERFORMANCE METRIC				
LAMBADA OPENAI	0.437	0.479	0.510	0.544
OPENBOOKQA	0.316	0.314	0.326	0.328
WINOGRANDE	0.515	0.533	0.530	0.576
PIQA	0.677	0.699	0.707	0.711
BOOLQ	0.547	0.492	0.584	0.574
WSC273	0.568	0.608	0.634	0.667
HELLASWAG 0SHOT	0.397	0.432	0.477	0.505
ARC CHALLENGE	0.267	0.280	0.292	0.308
ARC EASY	0.541	0.569	0.606	0.618
COPA	0.670	0.730	0.760	0.720
FINAL C4 LOSS	3.260	3.165	3.094	3.043

Table 26: Evaluation Performance for Cautious, Model Size = 520m

DATA SIZE	10B	21B	42B	85B
PERFORMANCE METRIC				
LAMBADA OPENAI	0.509	0.549	0.578	0.604
OPENBOOKQA	0.318	0.338	0.342	0.364
WINOGRANDE	0.548	0.575	0.571	0.605
PIQA	0.705	0.715	0.733	0.739
BOOLQ	0.601	0.611	0.582	0.606
WSC273	0.619	0.718	0.659	0.736
HELLASWAG 0SHOT	0.467	0.514	0.550	0.583
ARC CHALLENGE	0.305	0.315	0.344	0.358
ARC EASY	0.609	0.649	0.676	0.689
COPA	0.730	0.760	0.770	0.770
FINAL C4 LOSS	3.100	3.017	2.956	2.910

Table 27: Evaluation Performance for SOAP, Model Size = 130m

DATA SIZE	2B	5B	10B	21B
PERFORMANCE METRIC				
LAMBADA OPENAI	0.344	0.404	0.434	0.458
OPENBOOKQA	0.280	0.270	0.300	0.318
WINOGRANDE	0.511	0.510	0.523	0.517
PIQA	0.643	0.652	0.666	0.676
BOOLQ	0.569	0.491	0.519	0.559
WSC273	0.513	0.568	0.597	0.612
HELLASWAG 0SHOT	0.321	0.351	0.376	0.401
ARC CHALLENGE	0.247	0.255	0.270	0.270
ARC EASY	0.475	0.519	0.537	0.571
COPA	0.690	0.710	0.720	0.710
FINAL C4 LOSS	3.487	3.376	3.295	3.240

Table 28: Evaluation Performance for SOAP, Model Size = 300m

DATA SIZE	6B	12B	24B	48B
PERFORMANCE METRIC				
LAMBADA OPENAI	0.001	0.492	0.518	0.541
OPENBOOKQA	0.256	0.324	0.324	0.334
WINOGRANDE	0.506	0.534	0.546	0.555
PIQA	0.546	0.689	0.706	0.711
BOOLQ	0.402	0.593	0.565	0.593
WSC273	0.498	0.619	0.703	0.670
HELLASWAG 0SHOT	0.257	0.443	0.478	0.506
ARC CHALLENGE	0.219	0.281	0.305	0.324
ARC EASY	0.309	0.590	0.612	0.632
COPA	0.590	0.740	0.720	0.700
FINAL C4 LOSS	5.437	3.147	3.082	3.030

Table 29: Evaluation Performance for SOAP, Model Size = 520m

DATA SIZE	10B	21B	42B	85B
PERFORMANCE METRIC				
LAMBADA OPENAI	0.513	0.555	0.582	0.606
OPENBOOKQA	0.308	0.312	0.368	0.350
WINOGRANDE	0.548	0.562	0.594	0.583
PIQA	0.695	0.714	0.729	0.739
BOOLQ	0.563	0.623	0.596	0.643
WSC273	0.648	0.703	0.722	0.755
HELLASWAG 0SHOT	0.478	0.523	0.553	0.586
ARC CHALLENGE	0.292	0.345	0.334	0.356
ARC EASY	0.610	0.656	0.655	0.678
COPA	0.740	0.760	0.740	0.770
FINAL C4 LOSS	3.079	3.004	2.957	2.899

Table 30: Evaluation Performance for AdamW, Model Size = 130m

DATA SIZE	2B	5B	10B	21B
PERFORMANCE METRIC				
LAMBADA OPENAI	0.322	0.380	0.410	0.442
OPENBOOKQA	0.306	0.290	0.298	0.292
WINOGRANDE	0.506	0.519	0.515	0.530
PIQA	0.632	0.663	0.665	0.674
BOOLQ	0.568	0.415	0.553	0.576
WSC273	0.582	0.553	0.564	0.623
HELLASWAG 0SHOT	0.312	0.340	0.365	0.391
ARC CHALLENGE	0.229	0.255	0.262	0.270
ARC EASY	0.463	0.497	0.523	0.544
COPA	0.660	0.650	0.720	0.730
FINAL C4 LOSS	3.529	3.409	3.322	3.262

Table 31: Evaluation Performance for AdamW, Model Size = 300m

DATA SIZE	6B	12B	24B	48B
PERFORMANCE METRIC				
LAMBADA OPENAI	0.435	0.484	0.500	0.537
OPENBOOKQA	0.300	0.312	0.328	0.318
WINOGRANDE	0.519	0.530	0.542	0.552
PIQA	0.673	0.691	0.697	0.713
BOOLQ	0.454	0.496	0.508	0.602
WSC273	0.586	0.626	0.667	0.692
HELLASWAG 0SHOT	0.386	0.432	0.471	0.501
ARC CHALLENGE	0.269	0.295	0.304	0.310
ARC EASY	0.535	0.581	0.612	0.616
COPA	0.700	0.700	0.690	0.710
FINAL C4 LOSS	3.264	3.166	3.094	3.043

Table 32: Evaluation Performance for AdamW, Model Size = 520m

DATA SIZE	10B	21B	42B	85B
PERFORMANCE METRIC				
LAMBADA OPENAI	0.500	0.540	0.570	0.591
OPENBOOKQA	0.320	0.336	0.342	0.356
WINOGRANDE	0.539	0.586	0.583	0.601
PIQA	0.705	0.719	0.732	0.739
BOOLQ	0.555	0.615	0.596	0.609
WSC273	0.604	0.656	0.703	0.729
HELLASWAG 0SHOT	0.456	0.507	0.543	0.578
ARC CHALLENGE	0.299	0.311	0.323	0.349
ARC EASY	0.613	0.639	0.669	0.688
COPA	0.740	0.680	0.710	0.780
FINAL C4 LOSS	3.110	3.023	2.958	2.913

Table 33: Evaluation Performance for Adam-Mini, Model Size = 130m

DATA SIZE	2B	5B	10B	21B
PERFORMANCE METRIC				
LAMBADA OPENAI	0.317	0.369	0.412	0.444
OPENBOOKQA	0.286	0.282	0.282	0.302
WINOGRANDE	0.525	0.500	0.515	0.528
PIQA	0.640	0.655	0.662	0.670
BOOLQ	0.506	0.496	0.445	0.583
WSC273	0.571	0.557	0.546	0.634
HELLASWAG 0SHOT	0.310	0.337	0.364	0.390
ARC CHALLENGE	0.241	0.247	0.255	0.264
ARC EASY	0.465	0.495	0.525	0.555
COPA	0.690	0.710	0.730	0.700
FINAL C4 LOSS	3.542	3.416	3.328	3.266

Table 34: Evaluation Performance for Adam-Mini, Model Size = 300m

DATA SIZE	6B	12B	24B	48B
PERFORMANCE METRIC				
LAMBADA OPENAI	0.429	0.468	0.509	0.535
OPENBOOKQA	0.300	0.318	0.318	0.350
WINOGRANDE	0.502	0.528	0.545	0.586
PIQA	0.688	0.693	0.683	0.701
BOOLQ	0.480	0.546	0.609	0.573
WSC273	0.586	0.546	0.593	0.689
HELLASWAG 0SHOT	0.381	0.424	0.464	0.492
ARC CHALLENGE	0.270	0.292	0.294	0.300
ARC EASY	0.545	0.579	0.607	0.625
COPA	0.700	0.720	0.740	0.730
FINAL C4 LOSS	3.272	3.178	3.103	3.049

Table 35: Evaluation Performance for Adam-Mini, Model Size = 520m

DATA SIZE	10B	21B	42B	85B
PERFORMANCE METRIC				
LAMBADA OPENAI	0.500	0.538	0.576	0.604
OPENBOOKQA	0.306	0.326	0.330	0.368
WINOGRANDE	0.534	0.561	0.571	0.594
PIQA	0.707	0.721	0.723	0.733
BOOLQ	0.543	0.626	0.619	0.536
WSC273	0.619	0.700	0.733	0.744
HELLASWAG 0SHOT	0.459	0.505	0.541	0.576
ARC CHALLENGE	0.288	0.327	0.332	0.352
ARC EASY	0.613	0.638	0.654	0.671
COPA	0.710	0.720	0.760	0.760
FINAL C4 LOSS	3.112	3.027	2.966	2.912

## C Hyperparameter Ablation in Phase I

We reported the results for the optimizers we swept in Phase I. The result is formulated as follows: the first row shows the approximately best configuration found and the following rows show the results for the 1-dimensional ablations centered around the found configuration. The loss presented here is the final loss on the C4/EN validation set.



## C.1 Sweeping Results for AdamW

Table 36: Hyperparameter ablation for AdamW on 130m on 1x Chinchilla Data

$\beta_1$	$\beta_2$	$\epsilon$	$\eta$	$g_{\text{norm}}$	$\eta_{\text{min}}$	BSZ	warmup	$\lambda$	Loss	Link
0.9	0.98	1e-20	0.008	1	0	128	2000	0.1	3.529	0
0.95	–	–	–	–	–	–	–	–	3.539	1
0.98	–	–	–	–	–	–	–	–	3.882	2
–	0.9	–	–	–	–	–	–	–	3.545	3
–	0.95	–	–	–	–	–	–	–	3.535	4
–	–	1e-25	–	–	–	–	–	–	3.529	5
–	–	1e-15	–	–	–	–	–	–	3.531	6
–	–	1e-10	–	–	–	–	–	–	3.531	7
–	–	–	0.004	–	–	–	–	–	3.550	8
–	–	–	0.016	–	–	–	–	–	3.538	9
–	–	–	0.032	–	–	–	–	–	7.781	10
–	–	–	–	0	–	–	–	–	3.534	11
–	–	–	–	2.0	–	–	–	–	3.534	12
–	–	–	–	–	–	256	–	–	3.611	13
–	–	–	–	–	–	–	500	–	7.452	14
–	–	–	–	–	–	–	1000	–	3.532	15
–	–	–	–	–	–	–	4000	–	3.575	16
–	–	–	–	–	–	–	–	0	3.545	17
–	–	–	–	–	–	–	–	0.2	3.536	18

Table 37: Hyperparameter ablation for AdamW on 130m on 2x Chinchilla Data

$\beta_1$	$\beta_2$	$\epsilon$	$\eta$	$g_{\text{norm}}$	$\eta_{\text{min}}$	BSZ	warmup	$\lambda$	Loss	Link
0.9	0.98	1e-20	0.008	1	0	128	2000	0.1	3.409	0
0.95	–	–	–	–	–	–	–	–	3.417	1
0.98	–	–	–	–	–	–	–	–	7.557	2
–	0.9	–	–	–	–	–	–	–	3.423	3
–	0.95	–	–	–	–	–	–	–	3.413	4
–	–	1e-25	–	–	–	–	–	–	3.409	5
–	–	1e-15	–	–	–	–	–	–	3.409	6
–	–	1e-10	–	–	–	–	–	–	3.410	7
–	–	–	0.004	–	–	–	–	–	3.420	8
–	–	–	0.016	–	–	–	–	–	3.419	9
–	–	–	0.032	–	–	–	–	–	7.840	10
–	–	–	–	0	–	–	–	–	3.410	11
–	–	–	–	2.0	–	–	–	–	3.408	12
–	–	–	–	–	–	256	–	–	3.437	13
–	–	–	–	–	–	512	–	–	3.527	14
–	–	–	–	–	–	–	500	–	7.277	15
–	–	–	–	–	–	–	1000	–	3.413	16
–	–	–	–	–	–	–	4000	–	3.415	17
–	–	–	–	–	–	–	–	0	3.436	18
–	–	–	–	–	–	–	–	0.2	3.415	19

Table 38: Hyperparameter ablation for AdamW on 130m on 4x Chinchilla Data

$\beta_1$	$\beta_2$	$\epsilon$	$\eta$	$g_{\text{norm}}$	$\eta_{\text{min}}$	BSZ	warmup	$\lambda$	Loss	Link
0.9	0.98	1e-20	0.008	1	0	128	2000	0.1	3.322	0
0.95	–	–	–	–	–	–	–	–	3.330	1
0.98	–	–	–	–	–	–	–	–	3.416	2
–	0.9	–	–	–	–	–	–	–	3.338	3
–	0.95	–	–	–	–	–	–	–	3.329	4
–	–	1e-25	–	–	–	–	–	–	3.322	5
–	–	1e-15	–	–	–	–	–	–	3.323	6
–	–	1e-10	–	–	–	–	–	–	3.324	7
–	–	–	0.004	–	–	–	–	–	3.329	8
–	–	–	0.016	–	–	–	–	–	3.337	9
–	–	–	0.032	–	–	–	–	–	7.562	10
–	–	–	–	0	–	–	–	–	3.327	11
–	–	–	–	2.0	–	–	–	–	3.324	12
–	–	–	–	–	–	256	–	–	3.331	13
–	–	–	–	–	–	512	–	–	3.373	14
–	–	–	–	–	–	1024	–	–	3.480	15
–	–	–	–	–	–	–	500	–	7.262	16
–	–	–	–	–	–	–	1000	–	3.327	17
–	–	–	–	–	–	–	4000	–	3.325	18
–	–	–	–	–	–	–	–	0	3.359	19
–	–	–	–	–	–	–	–	0.2	3.335	20

Table 39: Hyperparameter ablation for AdamW on 130m on 8x Chinchilla Data

$\beta_1$	$\beta_2$	$\epsilon$	$\eta$	$g_{\text{norm}}$	$\eta_{\text{min}}$	BSZ	warmup	$\lambda$	Loss	Link
0.9	0.98	1e-20	0.008	1	0	256	1000	0.1	3.262	0
0.95	–	–	–	–	–	–	–	–	3.273	1
0.98	–	–	–	–	–	–	–	–	3.430	2
–	0.9	–	–	–	–	–	–	–	3.272	3
–	0.95	–	–	–	–	–	–	–	3.266	4
–	–	1e-25	–	–	–	–	–	–	3.262	5
–	–	1e-15	–	–	–	–	–	–	3.263	6
–	–	1e-10	–	–	–	–	–	–	3.261	7
–	–	–	0.004	–	–	–	–	–	3.270	8
–	–	–	0.016	–	–	–	–	–	7.435	9
–	–	–	0.032	–	–	–	–	–	7.658	10
–	–	–	–	0	–	–	–	–	3.263	11
–	–	–	–	2.0	–	–	–	–	3.264	12
–	–	–	–	–	–	128	–	–	3.264	13
–	–	–	–	–	–	512	–	–	3.286	14
–	–	–	–	–	–	1024	–	–	3.328	15
–	–	–	–	–	–	–	500	–	3.278	16
–	–	–	–	–	–	–	2000	–	3.263	17
–	–	–	–	–	–	–	4000	–	3.262	18
–	–	–	–	–	–	–	–	0	3.310	19
–	–	–	–	–	–	–	–	0.2	3.269	20

Table 40: Hyperparameter ablation for AdamW on 300m on 1x Chinchilla Data

$\beta_1$	$\beta_2$	$\epsilon$	$\eta$	$g_{\text{norm}}$	$\eta_{\text{min}}$	BSZ	warmup	$\lambda$	Loss	Link
0.9	0.98	1e-10	0.008	1	0	128	2000	0.1	3.264	0
0.95	–	–	–	–	–	–	–	–	3.271	1
0.98	–	–	–	–	–	–	–	–	7.351	2
–	0.9	–	–	–	–	–	–	–	3.280	3
–	0.95	–	–	–	–	–	–	–	3.269	4
–	–	1e-25	–	–	–	–	–	–	3.265	5
–	–	1e-20	–	–	–	–	–	–	3.265	6
–	–	1e-15	–	–	–	–	–	–	3.263	7
–	–	–	0.004	–	–	–	–	–	3.272	8
–	–	–	0.016	–	–	–	–	–	7.760	9
–	–	–	0.032	–	–	–	–	–	7.784	10
–	–	–	–	0	–	–	–	–	3.263	11
–	–	–	–	2.0	–	–	–	–	3.263	12
–	–	–	–	–	–	256	–	–	3.282	13
–	–	–	–	–	–	512	–	–	3.367	14
–	–	–	–	–	–	–	500	–	7.704	15
–	–	–	–	–	–	–	1000	–	7.759	16
–	–	–	–	–	–	–	4000	–	3.270	17
–	–	–	–	–	–	–	–	0	3.303	18
–	–	–	–	–	–	–	–	0.2	3.275	19

Table 41: Hyperparameter ablation for AdamW on 520m on 1x Chinchilla Data

$\beta_1$	$\beta_2$	$\epsilon$	$\eta$	$g_{\text{norm}}$	$\eta_{\text{min}}$	BSZ	warmup	$\lambda$	Loss	Link
0.9	0.98	1e-10	0.004	1	0	256	1000	0.2	3.110	0
0.95	–	–	–	–	–	–	–	–	3.112	1
0.98	–	–	–	–	–	–	–	–	3.229	2
–	0.9	–	–	–	–	–	–	–	3.116	3
–	0.95	–	–	–	–	–	–	–	3.111	4
–	–	1e-25	–	–	–	–	–	–	3.116	5
–	–	1e-20	–	–	–	–	–	–	3.116	6
–	–	1e-15	–	–	–	–	–	–	3.115	7
–	–	–	0.008	–	–	–	–	–	7.837	8
–	–	–	0.016	–	–	–	–	–	7.756	9
–	–	–	0.032	–	–	–	–	–	7.680	10
–	–	–	–	0	–	–	–	–	3.114	11
–	–	–	–	2.0	–	–	–	–	3.118	12
–	–	–	–	–	–	128	–	–	7.630	13
–	–	–	–	–	–	512	–	–	3.169	14
–	–	–	–	–	–	1024	–	–	3.302	15
–	–	–	–	–	–	–	500	–	3.165	16
–	–	–	–	–	–	–	2000	–	3.113	17
–	–	–	–	–	–	–	4000	–	3.126	18
–	–	–	–	–	–	–	–	0	7.270	19
–	–	–	–	–	–	–	–	0.1	3.135	20

## C.2 Sweeping Results for Cautious

Table 42: Hyperparameter ablation for Cautious on 130m on 1x Chinchilla Data

$\beta_1$	$\beta_2$	$\epsilon$	$\eta$	$g_{\text{norm}}$	$\eta_{\text{min}}$	BSZ	warmup	$\lambda$	Loss	Link
0.95	0.98	1e-15	0.008	1	0	128	2000	0.1	3.535	0
0.8	–	–	–	–	–	–	–	–	6.698	1
0.9	–	–	–	–	–	–	–	–	3.549	2
0.98	–	–	–	–	–	–	–	–	3.534	3
–	0.9	–	–	–	–	–	–	–	3.551	4
–	0.95	–	–	–	–	–	–	–	3.543	5
–	–	1e-25	–	–	–	–	–	–	3.537	6
–	–	1e-20	–	–	–	–	–	–	3.537	7
–	–	1e-10	–	–	–	–	–	–	3.536	8
–	–	–	0.016	–	–	–	–	–	3.539	9
–	–	–	0.032	–	–	–	–	–	7.802	10
–	–	–	–	0	–	–	–	–	3.534	11
–	–	–	–	2.0	–	–	–	–	3.532	12
–	–	–	–	–	–	256	–	–	3.610	13
–	–	–	–	–	–	–	500	–	3.572	14
–	–	–	–	–	–	–	1000	–	3.535	15
–	–	–	–	–	–	–	4000	–	3.582	16
–	–	–	–	–	–	–	–	0	3.552	17
–	–	–	–	–	–	–	–	0.2	3.537	18

Table 43: Hyperparameter ablation for Cautious on 130m on 2x Chinchilla Data

$\beta_1$	$\beta_2$	$\epsilon$	$\eta$	$g_{\text{norm}}$	$\eta_{\text{min}}$	BSZ	warmup	$\lambda$	Loss	Link
0.98	0.98	1e-15	0.008	2	0	128	2000	0.1	3.403	0
0.8	–	–	–	–	–	–	–	–	7.417	1
0.9	–	–	–	–	–	–	–	–	3.427	2
0.95	–	–	–	–	–	–	–	–	3.413	3
–	0.9	–	–	–	–	–	–	–	>10	4
–	0.95	–	–	–	–	–	–	–	3.408	5
–	–	1e-25	–	–	–	–	–	–	3.404	6
–	–	1e-20	–	–	–	–	–	–	3.404	7
–	–	1e-10	–	–	–	–	–	–	3.404	8
–	–	–	0.016	–	–	–	–	–	3.416	9
–	–	–	0.032	–	–	–	–	–	3.513	10
–	–	–	–	0	–	–	–	–	3.415	11
–	–	–	–	1.0	–	–	–	–	3.401	12
–	–	–	–	–	–	256	–	–	3.422	13
–	–	–	–	–	–	512	–	–	3.499	14
–	–	–	–	–	–	–	500	–	3.453	15
–	–	–	–	–	–	–	1000	–	3.412	16
–	–	–	–	–	–	–	4000	–	3.410	17
–	–	–	–	–	–	–	–	0	3.436	18
–	–	–	–	–	–	–	–	0.2	3.405	19

Table 44: Hyperparameter ablation for Cautious on 130m on 8x Chinchilla Data

$\beta_1$	$\beta_2$	$\epsilon$	$\eta$	$g_{\text{norm}}$	$\eta_{\text{min}}$	BSZ	warmup	$\lambda$	Loss	Link
0.98	0.98	1e-15	0.008	2	0	256	2000	0.1	3.253	0
0.8	–	–	–	–	–	–	–	–	7.395	1
0.9	–	–	–	–	–	–	–	–	3.271	2
0.95	–	–	–	–	–	–	–	–	3.259	3
–	0.9	–	–	–	–	–	–	–	>10	4
–	0.95	–	–	–	–	–	–	–	3.253	5
–	–	1e-25	–	–	–	–	–	–	3.251	6
–	–	1e-20	–	–	–	–	–	–	3.251	7
–	–	1e-10	–	–	–	–	–	–	3.250	8
–	–	–	0.016	–	–	–	–	–	3.257	9
–	–	–	0.032	–	–	–	–	–	3.383	10
–	–	–	–	0	–	–	–	–	3.254	11
–	–	–	–	1.0	–	–	–	–	3.251	12
–	–	–	–	–	–	128	–	–	3.265	13
–	–	–	–	–	–	512	–	–	3.261	14
–	–	–	–	–	–	1024	–	–	3.294	15
–	–	–	–	–	–	–	500	–	3.276	16
–	–	–	–	–	–	–	1000	–	3.253	17
–	–	–	–	–	–	–	4000	–	3.255	18
–	–	–	–	–	–	–	–	0	3.291	19
–	–	–	–	–	–	–	–	0.2	3.253	20

Table 45: Hyperparameter ablation for Cautious on 300m on 1x Chinchilla Data

$\beta_1$	$\beta_2$	$\epsilon$	$\eta$	$g_{\text{norm}}$	$\eta_{\text{min}}$	BSZ	warmup	$\lambda$	Loss	Link
0.98	0.98	1e-25	0.008	2	0	128	2000	0.1	3.260	0
0.9	–	–	–	–	–	–	–	–	3.286	1
0.95	–	–	–	–	–	–	–	–	3.271	2
–	0.9	–	–	–	–	–	–	–	>10	3
–	0.95	–	–	–	–	–	–	–	>10	4
–	–	1e-25	–	–	–	–	–	–	3.260	5
–	–	1e-15	–	–	–	–	–	–	3.260	6
–	–	–	–	0	–	–	–	–	3.274	7
–	–	–	–	1	–	–	–	–	3.259	8
–	–	–	–	–	–	256	–	–	3.270	9
–	–	–	–	–	–	–	1000	–	7.352	10
–	–	–	–	–	–	–	4000	–	3.264	11
–	–	–	–	–	–	–	–	0.0	3.310	12
–	–	–	–	–	–	–	–	0.2	3.266	13

Table 46: Hyperparameter ablation for Cautious on 520m on 1x Chinchilla Data

$\beta_1$	$\beta_2$	$\epsilon$	$\eta$	$g_{\text{norm}}$	$\eta_{\text{min}}$	BSZ	warmup	$\lambda$	Loss	Link
0.98	0.98	1e-25	0.008	1	0	256	2000	0.1	3.100	0
0.8	–	–	–	–	–	–	–	–	>10	1
0.9	–	–	–	–	–	–	–	–	7.264	2
0.95	–	–	–	–	–	–	–	–	3.108	3
–	0.9	–	–	–	–	–	–	–	>10	4
–	0.95	–	–	–	–	–	–	–	3.105	5
–	–	1e-25	–	–	–	–	–	–	3.100	6
–	–	1e-20	–	–	–	–	–	–	3.100	7
–	–	1e-15	–	–	–	–	–	–	3.101	8
–	–	1e-10	–	–	–	–	–	–	3.101	9
–	–	–	0.016	–	–	–	–	–	3.125	10
–	–	–	0.032	–	–	–	–	–	7.662	11
–	–	–	–	0	–	–	–	–	3.123	12
–	–	–	–	2.0	–	–	–	–	3.102	13
–	–	–	–	–	–	128	–	–	>10	14
–	–	–	–	–	–	512	–	–	3.123	15
–	–	–	–	–	–	–	500	–	>10	16
–	–	–	–	–	–	–	1000	–	3.119	17
–	–	–	–	–	–	–	4000	–	3.107	18
–	–	–	–	–	–	–	–	0	3.131	19
–	–	–	–	–	–	–	–	0.2	3.105	20

### C.3 Sweeping Results for Lion

Table 47: Hyperparameter ablation for Lion on 130m on 1x Chinchilla Data

$\beta_1$	$\beta_2$	$\eta$	$g_{\text{norm}}$	$\eta_{\text{min}}$	BSZ	warmup	$\lambda$	Loss	Link
0.9	0.95	0.002	1	0	128	2000	0.7	3.552	0
0.8	—	—	—	—	—	—	—	3.575	1
0.95	—	—	—	—	—	—	—	3.557	2
0.98	—	—	—	—	—	—	—	7.644	3
—	0.9	—	—	—	—	—	—	3.550	4
—	0.98	—	—	—	—	—	—	3.630	5
—	—	0.0005	—	—	—	—	—	3.579	6
—	—	0.001	—	—	—	—	—	3.549	7
—	—	0.004	—	—	—	—	—	7.806	8
—	—	0.008	—	—	—	—	—	7.828	9
—	—	—	0	—	—	—	—	3.559	10
—	—	—	2.0	—	—	—	—	3.557	11
—	—	—	—	—	256	—	—	3.643	12
—	—	—	—	—	—	500	—	7.840	13
—	—	—	—	—	—	1000	—	7.739	14
—	—	—	—	—	—	4000	—	3.601	15
—	—	—	—	—	—	—	0	3.570	16
—	—	—	—	—	—	—	0.1	3.562	17
—	—	—	—	—	—	—	0.2	3.557	18
—	—	—	—	—	—	—	0.3	3.555	19
—	—	—	—	—	—	—	0.4	3.553	20
—	—	—	—	—	—	—	0.5	3.551	21
—	—	—	—	—	—	—	0.6	3.549	22
—	—	—	—	—	—	—	0.8	3.554	23
—	—	—	—	—	—	—	0.9	3.556	24
—	—	—	—	—	—	—	1	3.557	25

Table 48: Hyperparameter ablation for Lion on 130m on 2x Chinchilla Data

$\beta_1$	$\beta_2$	$\eta$	$g_{\text{norm}}$	$\eta_{\text{min}}$	BSZ	warmup	$\lambda$	Loss	Link
0.9	0.95	0.001	1	0	128	2000	0.7	3.414	0
0.8	–	–	–	–	–	–	–	3.436	1
0.95	–	–	–	–	–	–	–	3.418	2
0.98	–	–	–	–	–	–	–	7.313	3
–	0.9	–	–	–	–	–	–	3.433	4
–	0.98	–	–	–	–	–	–	3.409	5
–	–	0.0005	–	–	–	–	–	3.435	6
–	–	0.002	–	–	–	–	–	3.414	7
–	–	0.004	–	–	–	–	–	3.489	8
–	–	0.008	–	–	–	–	–	7.826	9
–	–	–	0	–	–	–	–	3.414	10
–	–	–	2.0	–	–	–	–	3.413	11
–	–	–	–	–	256	–	–	3.447	12
–	–	–	–	–	512	–	–	3.540	13
–	–	–	–	–	–	500	–	3.435	14
–	–	–	–	–	–	1000	–	3.415	15
–	–	–	–	–	–	4000	–	3.423	16
–	–	–	–	–	–	–	0.4	3.419	17
–	–	–	–	–	–	–	0.5	3.418	18
–	–	–	–	–	–	–	0.6	3.414	19
–	–	–	–	–	–	–	0.8	3.414	20
–	–	–	–	–	–	–	0.9	3.413	21
–	–	–	–	–	–	–	1	3.413	22



Table 49: Hyperparameter ablation for Lion on 130m on 4x Chinchilla Data

$\beta_1$	$\beta_2$	$\eta$	$g_{\text{norm}}$	$\eta_{\text{min}}$	BSZ	warmup	$\lambda$	Loss	Link
0.9	0.95	0.001	1	0	128	2000	0.7	3.331	0
0.8	–	–	–	–	–	–	–	3.345	1
0.95	–	–	–	–	–	–	–	3.335	2
0.98	–	–	–	–	–	–	–	7.330	3
–	0.9	–	–	–	–	–	–	3.345	4
–	0.98	–	–	–	–	–	–	3.340	5
–	–	0.0005	–	–	–	–	–	3.338	6
–	–	0.002	–	–	–	–	–	3.342	7
–	–	0.004	–	–	–	–	–	7.406	8
–	–	0.008	–	–	–	–	–	>10	9
–	–	–	0	–	–	–	–	3.333	10
–	–	–	2.0	–	–	–	–	3.334	11
–	–	–	–	–	256	–	–	3.346	12
–	–	–	–	–	512	–	–	3.386	13
–	–	–	–	–	1024	–	–	3.492	14
–	–	–	–	–	–	500	–	3.364	15
–	–	–	–	–	–	1000	–	3.336	16
–	–	–	–	–	–	4000	–	3.333	17
–	–	–	–	–	–	–	0.4	3.335	18
–	–	–	–	–	–	–	0.5	3.333	19
–	–	–	–	–	–	–	0.6	3.329	20
–	–	–	–	–	–	–	0.8	3.332	21
–	–	–	–	–	–	–	0.9	3.332	22
–	–	–	–	–	–	–	1	3.335	23

Table 50: Hyperparameter ablation for Lion on 130m on 8x Chinchilla Data

$\beta_1$	$\beta_2$	$\eta$	$g_{\text{norm}}$	$\eta_{\text{min}}$	BSZ	warmup	$\lambda$	Loss	Link
0.9	0.98	0.001	1	0	128	2000	0.7	3.252	0
0.8	–	–	–	–	–	–	–	3.287	1
0.95	–	–	–	–	–	–	–	3.264	2
0.98	–	–	–	–	–	–	–	3.286	3
–	0.9	–	–	–	–	–	–	3.277	4
–	0.95	–	–	–	–	–	–	3.263	5
–	–	0.0005	–	–	–	–	–	3.254	6
–	–	0.002	–	–	–	–	–	3.310	7
–	–	0.004	–	–	–	–	–	7.829	8
–	–	0.008	–	–	–	–	–	NaN	9
–	–	–	0	–	–	–	–	3.329	10
–	–	–	2.0	–	–	–	–	3.265	11
–	–	–	–	–	256	–	–	3.260	12
–	–	–	–	–	512	–	–	3.287	13
–	–	–	–	–	1024	–	–	3.342	14
–	–	–	–	–	–	500	–	3.336	15
–	–	–	–	–	–	1000	–	3.273	16
–	–	–	–	–	–	4000	–	3.258	17
–	–	–	–	–	–	–	0.4	3.256	18
–	–	–	–	–	–	–	0.5	3.251	19
–	–	–	–	–	–	–	0.6	3.252	20
–	–	–	–	–	–	–	0.8	3.258	21
–	–	–	–	–	–	–	0.9	3.259	22
–	–	–	–	–	–	–	1	3.261	23

Table 51: Hyperparameter ablation for Lion on 300m on 1x Chinchilla Data

$\beta_1$	$\beta_2$	$\eta$	$g_{\text{norm}}$	$\eta_{\text{min}}$	BSZ	warmup	$\lambda$	Loss	Link
0.9	0.95	0.001	1	0	128	2000	0.6	3.268	0
0.8	–	–	–	–	–	–	–	3.283	1
0.95	–	–	–	–	–	–	–	3.271	2
0.98	–	–	–	–	–	–	–	7.690	3
–	0.9	–	–	–	–	–	–	3.283	4
–	0.98	–	–	–	–	–	–	3.271	5
–	–	0.0005	–	–	–	–	–	3.286	6
–	–	0.002	–	–	–	–	–	3.283	7
–	–	0.004	–	–	–	–	–	7.817	8
–	–	0.008	–	–	–	–	–	7.863	9
–	–	–	0	–	–	–	–	3.274	10
–	–	–	2.0	–	–	–	–	3.269	11
–	–	–	–	–	256	–	–	3.295	12
–	–	–	–	–	512	–	–	3.367	13
–	–	–	–	–	–	500	–	7.607	14
–	–	–	–	–	–	1000	–	3.278	15
–	–	–	–	–	–	4000	–	3.277	16
–	–	–	–	–	–	–	0.4	3.272	17
–	–	–	–	–	–	–	0.5	3.271	18
–	–	–	–	–	–	–	0.7	3.268	19
–	–	–	–	–	–	–	0.8	3.268	20
–	–	–	–	–	–	–	0.9	3.269	21
–	–	–	–	–	–	–	1	3.269	22

Table 52: Hyperparameter ablation for Lion on 520m on 1x Chinchilla Data

$\beta_1$	$\beta_2$	$\eta$	$g_{\text{norm}}$	$\eta_{\text{min}}$	BSZ	warmup	$\lambda$	Loss	Link
0.9	0.95	0.001	1	0	128	2000	0.7	3.108	0
0.8	–	–	–	–	–	–	–	3.117	1
0.95	–	–	–	–	–	–	–	3.121	2
0.98	–	–	–	–	–	–	–	7.577	3
–	0.9	–	–	–	–	–	–	3.119	4
–	0.98	–	–	–	–	–	–	3.208	5
–	–	0.0005	–	–	–	–	–	3.113	6
–	–	0.002	–	–	–	–	–	7.734	7
–	–	0.004	–	–	–	–	–	8.046	8
–	–	0.008	–	–	–	–	–	NaN	9
–	–	–	0	–	–	–	–	3.128	10
–	–	–	2.0	–	–	–	–	3.109	11
–	–	–	–	–	256	–	–	3.117	12
–	–	–	–	–	512	–	–	3.151	13
–	–	–	–	–	1024	–	–	3.252	14
–	–	–	–	–	–	500	–	7.403	15
–	–	–	–	–	–	1000	–	7.119	16
–	–	–	–	–	–	4000	–	3.115	17
–	–	–	–	–	–	–	0.4	3.109	18
–	–	–	–	–	–	–	0.5	3.108	19
–	–	–	–	–	–	–	0.6	3.108	20
–	–	–	–	–	–	–	0.8	3.109	21
–	–	–	–	–	–	–	0.9	3.112	22
–	–	–	–	–	–	–	1	3.115	23

## C.4 Sweeping Results for Mars

Table 53: Hyperparameter ablation for Mars on 130m on 1x Chinchilla Data

$\beta_1$	$\beta_2$	$\epsilon$	$\gamma$	$\eta$	$g_{\text{norm}}$	$\eta_{\text{min}}$	BSZ	warmup	$\lambda$	Loss	Link
0.9	0.95	1e-25	0.025	0.016	1	0	128	2000	0.1	3.537	0
0.8	—	—	—	—	—	—	—	—	—	3.568	1
0.95	—	—	—	—	—	—	—	—	—	3.548	2
0.98	—	—	—	—	—	—	—	—	—	3.586	3
—	0.9	—	—	—	—	—	—	—	—	3.548	4
—	0.98	—	—	—	—	—	—	—	—	3.536	5
—	0.99	—	—	—	—	—	—	—	—	3.562	6
—	—	1e-30	—	—	—	—	—	—	—	3.537	7
—	—	1e-25	—	—	—	—	—	—	—	3.537	8
—	—	1e-20	—	—	—	—	—	—	—	3.537	9
—	—	1e-15	—	—	—	—	—	—	—	3.537	10
—	—	1e-10	—	—	—	—	—	—	—	3.537	11
—	—	—	0.0125	—	—	—	—	—	—	3.540	12
—	—	—	0.05	—	—	—	—	—	—	3.542	13
—	—	—	0.1	—	—	—	—	—	—	3.553	14
—	—	—	—	0.008	—	—	—	—	—	3.538	15
—	—	—	—	0.032	—	—	—	—	—	7.574	16
—	—	—	—	—	—	—	256	—	—	3.616	17
—	—	—	—	—	—	—	—	500	—	>10	18
—	—	—	—	—	—	—	—	1000	—	>10	19
—	—	—	—	—	—	—	—	4000	—	3.596	20
—	—	—	—	—	—	—	—	—	0	3.552	21
—	—	—	—	—	—	—	—	—	0.2	3.561	22

Table 54: Hyperparameter ablation for Mars on 130m on 2x Chinchilla Data

$\beta_1$	$\beta_2$	$\epsilon$	$\gamma$	$\eta$	$g_{\text{norm}}$	$\eta_{\text{min}}$	BSZ	warmup	$\lambda$	Loss	Link
0.95	0.98	1e-25	0.025	0.008	1	0	128	2000	0.1	3.396	0
0.8	–	–	–	–	–	–	–	–	–	3.421	1
0.9	–	–	–	–	–	–	–	–	–	3.402	2
0.98	–	–	–	–	–	–	–	–	–	3.401	3
–	0.9	–	–	–	–	–	–	–	–	3.409	4
–	0.95	–	–	–	–	–	–	–	–	3.400	5
–	0.99	–	–	–	–	–	–	–	–	3.398	6
–	–	1e-30	–	–	–	–	–	–	–	3.396	7
–	–	1e-20	–	–	–	–	–	–	–	3.396	8
–	–	1e-15	–	–	–	–	–	–	–	3.398	9
–	–	1e-10	–	–	–	–	–	–	–	3.397	10
–	–	–	0.0125	–	–	–	–	–	–	3.398	11
–	–	–	0.05	–	–	–	–	–	–	3.404	12
–	–	–	0.1	–	–	–	–	–	–	3.415	13
–	–	–	–	0.016	–	–	–	–	–	3.402	14
–	–	–	–	0.032	–	–	–	–	–	3.441	15
–	–	–	–	–	–	–	256	–	–	3.427	16
–	–	–	–	–	–	–	512	–	–	3.524	17
–	–	–	–	–	–	–	–	500	–	3.400	18
–	–	–	–	–	–	–	–	1000	–	3.395	19
–	–	–	–	–	–	–	–	4000	–	3.409	20
–	–	–	–	–	–	–	–	–	0	3.430	21
–	–	–	–	–	–	–	–	–	0.2	3.404	22

Table 55: Hyperparameter ablation for Mars on 130m on 4x Chinchilla Data

$\beta_1$	$\beta_2$	$\epsilon$	$\gamma$	$\eta$	$g_{\text{norm}}$	$\eta_{\text{min}}$	BSZ	warmup	$\lambda$	Loss	Link
0.9	0.98	1e-25	0.025	0.008	1	0	128	2000	0.1	3.323	0
0.8	–	–	–	–	–	–	–	–	–	3.336	1
0.95	–	–	–	–	–	–	–	–	–	3.326	2
0.98	–	–	–	–	–	–	–	–	–	3.350	3
–	0.9	–	–	–	–	–	–	–	–	3.337	4
–	0.95	–	–	–	–	–	–	–	–	3.330	5
–	0.99	–	–	–	–	–	–	–	–	3.324	6
–	–	1e-30	–	–	–	–	–	–	–	3.323	7
–	–	1e-25	–	–	–	–	–	–	–	3.323	8
–	–	1e-20	–	–	–	–	–	–	–	3.323	9
–	–	1e-15	–	–	–	–	–	–	–	3.321	10
–	–	1e-10	–	–	–	–	–	–	–	3.322	11
–	–	–	0.0125	–	–	–	–	–	–	3.323	12
–	–	–	0.05	–	–	–	–	–	–	3.324	13
–	–	–	0.1	–	–	–	–	–	–	3.330	14
–	–	–	–	0.016	–	–	–	–	–	3.337	15
–	–	–	–	0.032	–	–	–	–	–	8.642	16
–	–	–	–	–	–	–	256	–	–	3.336	17
–	–	–	–	–	–	–	512	–	–	3.384	18
–	–	–	–	–	–	–	1024	–	–	3.525	19
–	–	–	–	–	–	–	–	500	–	3.328	20
–	–	–	–	–	–	–	–	1000	–	3.321	21
–	–	–	–	–	–	–	–	4000	–	3.327	22
–	–	–	–	–	–	–	–	–	0	3.359	23
–	–	–	–	–	–	–	–	–	0.2	3.333	24

Table 56: Hyperparameter ablation for Mars on 130m on 8x Chinchilla Data

$\beta_1$	$\beta_2$	$\epsilon$	$\gamma$	$\eta$	$g_{\text{norm}}$	$\eta_{\text{min}}$	BSZ	warmup	$\lambda$	Loss	Link
0.98	0.99	1e-25	0.025	0.008	1	0	128	2000	0.1	3.247	0
0.8	–	–	–	–	–	–	–	–	–	3.272	1
0.9	–	–	–	–	–	–	–	–	–	3.252	2
0.95	–	–	–	–	–	–	–	–	–	3.256	3
–	0.9	–	–	–	–	–	–	–	–	NaN	4
–	0.95	–	–	–	–	–	–	–	–	3.257	5
–	0.98	–	–	–	–	–	–	–	–	3.249	6
–	–	1e-30	–	–	–	–	–	–	–	3.247	7
–	–	1e-25	–	–	–	–	–	–	–	3.247	8
–	–	1e-20	–	–	–	–	–	–	–	3.247	9
–	–	1e-15	–	–	–	–	–	–	–	3.247	10
–	–	1e-10	–	–	–	–	–	–	–	3.248	11
–	–	–	0.0125	–	–	–	–	–	–	3.255	12
–	–	–	0.05	–	–	–	–	–	–	3.252	13
–	–	–	0.1	–	–	–	–	–	–	3.410	14
–	–	–	–	0.016	–	–	–	–	–	3.262	15
–	–	–	–	0.032	–	–	–	–	–	3.322	16
–	–	–	–	–	–	–	256	–	–	3.250	17
–	–	–	–	–	–	–	512	–	–	3.277	18
–	–	–	–	–	–	–	1024	–	–	3.339	19
–	–	–	–	–	–	–	–	500	–	3.250	20
–	–	–	–	–	–	–	–	1000	–	3.248	21
–	–	–	–	–	–	–	–	4000	–	3.251	22
–	–	–	–	–	–	–	–	–	0	3.292	23
–	–	–	–	–	–	–	–	–	0.2	3.265	24



Table 57: Hyperparameter ablation for Mars on 300m on 1x Chinchilla Data

$\beta_1$	$\beta_2$	$\epsilon$	$\gamma$	$\eta$	$g_{\text{norm}}$	$\eta_{\text{min}}$	BSZ	warmup	$\lambda$	Loss	Link
0.98	0.98	1e-25	0.05	0.008	1	0	128	1000	0.1	3.249	0
0.8	–	–	–	–	–	–	–	–	–	7.899	1
0.9	–	–	–	–	–	–	–	–	–	3.259	2
0.95	–	–	–	–	–	–	–	–	–	3.247	3
–	0.9	–	–	–	–	–	–	–	–	NaN	4
–	0.95	–	–	–	–	–	–	–	–	3.256	5
–	0.99	–	–	–	–	–	–	–	–	3.247	6
–	–	1e-30	–	–	–	–	–	–	–	3.249	7
–	–	1e-20	–	–	–	–	–	–	–	3.250	8
–	–	1e-15	–	–	–	–	–	–	–	3.250	9
–	–	1e-10	–	–	–	–	–	–	–	3.252	10
–	–	–	0.0125	–	–	–	–	–	–	3.259	11
–	–	–	0.025	–	–	–	–	–	–	3.249	12
–	–	–	0.1	–	–	–	–	–	–	3.280	13
–	–	–	–	0.016	–	–	–	–	–	3.265	14
–	–	–	–	0.032	–	–	–	–	–	3.410	15
–	–	–	–	–	–	–	256	–	–	3.288	16
–	–	–	–	–	–	–	512	–	–	3.389	17
–	–	–	–	–	–	–	1024	–	–	3.624	18
–	–	–	–	–	–	–	–	500	–	3.254	19
–	–	–	–	–	–	–	–	2000	–	3.254	20
–	–	–	–	–	–	–	–	4000	–	3.269	21
–	–	–	–	–	–	–	–	–	0	3.312	22
–	–	–	–	–	–	–	–	–	0.2	3.262	23

Table 58: Hyperparameter ablation for Mars on 520m on 1x Chinchilla Data

$\beta_1$	$\beta_2$	$\epsilon$	$\gamma$	$\eta$	$g_{\text{norm}}$	$\eta_{\text{min}}$	BSZ	warmup	$\lambda$	Loss	Link
0.95	0.95	1e-25	0.025	0.008	1	0	256	2000	0.1	3.101	0
0.8	–	–	–	–	–	–	–	–	–	>10	1
0.9	–	–	–	–	–	–	–	–	–	3.108	2
0.98	–	–	–	–	–	–	–	–	–	3.107	3
–	0.9	–	–	–	–	–	–	–	–	3.108	4
–	0.98	–	–	–	–	–	–	–	–	3.099	5
–	0.99	–	–	–	–	–	–	–	–	3.102	6
–	–	1e-30	–	–	–	–	–	–	–	3.101	7
–	–	1e-20	–	–	–	–	–	–	–	3.101	8
–	–	1e-15	–	–	–	–	–	–	–	3.101	9
–	–	1e-10	–	–	–	–	–	–	–	3.101	10
–	–	–	0.0125	–	–	–	–	–	–	3.100	11
–	–	–	0.05	–	–	–	–	–	–	3.103	12
–	–	–	0.1	–	–	–	–	–	–	3.113	13
–	–	–	–	0.016	–	–	–	–	–	3.111	14
–	–	–	–	0.032	–	–	–	–	–	NaN	15
–	–	–	–	–	–	–	128	–	–	3.102	16
–	–	–	–	–	–	–	512	–	–	3.130	17
–	–	–	–	–	–	–	1024	–	–	3.224	18
–	–	–	–	–	–	–	–	500	–	3.128	19
–	–	–	–	–	–	–	–	1000	–	3.105	20
–	–	–	–	–	–	–	–	4000	–	3.112	21
–	–	–	–	–	–	–	–	–	0	3.131	22
–	–	–	–	–	–	–	–	–	0.2	3.103	23

## C.5 Sweeping Results for NAdamW

Table 59: Hyperparameter ablation for NAdamW on 130m on 1x Chinchilla Data

$\beta_1$	$\beta_2$	$\epsilon$	$\eta$	$g_{\text{norm}}$	$\eta_{\text{min}}$	BSZ	warmup	$\lambda$	Loss	Link
0.95	0.98	1e-25	0.008	1	0	128	2000	0.1	3.531	0
0.8	–	–	–	–	–	–	–	–	4.764	1
0.9	–	–	–	–	–	–	–	–	3.552	2
0.98	–	–	–	–	–	–	–	–	3.585	3
–	0.9	–	–	–	–	–	–	–	3.552	4
–	0.95	–	–	–	–	–	–	–	3.535	5
–	–	1e-20	–	–	–	–	–	–	3.531	6
–	–	1e-15	–	–	–	–	–	–	3.533	7
–	–	1e-10	–	–	–	–	–	–	3.531	8
–	–	–	0.016	–	–	–	–	–	3.545	9
–	–	–	0.032	–	–	–	–	–	>10	10
–	–	–	–	0	–	–	–	–	3.537	11
–	–	–	–	2.0	–	–	–	–	3.539	12
–	–	–	–	–	–	256	–	–	3.624	13
–	–	–	–	–	–	–	500	–	3.646	14
–	–	–	–	–	–	–	1000	–	3.545	15
–	–	–	–	–	–	–	4000	–	3.577	16
–	–	–	–	–	–	–	–	0	3.547	17
–	–	–	–	–	–	–	–	0.2	3.535	18

Table 60: Hyperparameter ablation for NAdamW on 130m on 2x Chinchilla Data

$\beta_1$	$\beta_2$	$\epsilon$	$\eta$	$g_{\text{norm}}$	$\eta_{\text{min}}$	BSZ	warmup	$\lambda$	Loss	Link
0.95	0.98	1e-10	0.008	1	0	128	2000	0.1	3.394	0
0.8	–	–	–	–	–	–	–	–	3.452	1
0.9	–	–	–	–	–	–	–	–	3.408	2
0.98	–	–	–	–	–	–	–	–	3.402	3
–	0.9	–	–	–	–	–	–	–	3.403	4
–	0.95	–	–	–	–	–	–	–	3.399	5
–	–	1e-25	–	–	–	–	–	–	3.394	6
–	–	1e-20	–	–	–	–	–	–	3.394	7
–	–	1e-15	–	–	–	–	–	–	3.393	8
–	–	–	0.016	–	–	–	–	–	3.406	9
–	–	–	0.032	–	–	–	–	–	7.675	10
–	–	–	–	0	–	–	–	–	3.400	11
–	–	–	–	2.0	–	–	–	–	3.396	12
–	–	–	–	–	–	256	–	–	3.423	13
–	–	–	–	–	–	512	–	–	3.520	14
–	–	–	–	–	–	–	500	–	3.424	15
–	–	–	–	–	–	–	1000	–	3.400	16
–	–	–	–	–	–	–	4000	–	3.405	17
–	–	–	–	–	–	–	–	0	3.421	18
–	–	–	–	–	–	–	–	0.2	3.398	19

Table 61: Hyperparameter ablation for NAdamW on 130m on 4x Chinchilla Data

$\beta_1$	$\beta_2$	$\epsilon$	$\eta$	$g_{\text{norm}}$	$\eta_{\text{min}}$	BSZ	warmup	$\lambda$	Loss	Link
0.95	0.98	1e-10	0.008	1	0	128	2000	0.1	3.319	0
0.8	–	–	–	–	–	–	–	–	7.085	1
0.9	–	–	–	–	–	–	–	–	3.331	2
0.98	–	–	–	–	–	–	–	–	3.349	3
–	0.9	–	–	–	–	–	–	–	3.332	4
–	0.95	–	–	–	–	–	–	–	3.327	5
–	–	1e-25	–	–	–	–	–	–	3.321	6
–	–	1e-20	–	–	–	–	–	–	3.321	7
–	–	1e-15	–	–	–	–	–	–	3.323	8
–	–	–	0.016	–	–	–	–	–	3.343	9
–	–	–	0.032	–	–	–	–	–	7.733	10
–	–	–	–	0	–	–	–	–	3.324	11
–	–	–	–	2.0	–	–	–	–	3.323	12
–	–	–	–	–	–	256	–	–	3.332	13
–	–	–	–	–	–	512	–	–	3.372	14
–	–	–	–	–	–	1024	–	–	3.496	15
–	–	–	–	–	–	–	500	–	6.917	16
–	–	–	–	–	–	–	1000	–	3.328	17
–	–	–	–	–	–	–	4000	–	3.321	18
–	–	–	–	–	–	–	–	0	3.359	19
–	–	–	–	–	–	–	–	0.2	3.330	20

Table 62: Hyperparameter ablation for NAdamW on 130m on 8x Chinchilla Data

$\beta_1$	$\beta_2$	$\epsilon$	$\eta$	$g_{\text{norm}}$	$\eta_{\text{min}}$	BSZ	warmup	$\lambda$	Loss	Link
0.95	0.98	1e-10	0.008	1	0	128	2000	0.1	3.251	0
0.8	–	–	–	–	–	–	–	–	3.282	1
0.9	–	–	–	–	–	–	–	–	3.257	2
0.98	–	–	–	–	–	–	–	–	3.253	3
–	0.9	–	–	–	–	–	–	–	3.260	4
–	0.95	–	–	–	–	–	–	–	3.255	5
–	–	1e-25	–	–	–	–	–	–	3.251	6
–	–	1e-20	–	–	–	–	–	–	3.251	7
–	–	1e-15	–	–	–	–	–	–	3.250	8
–	–	–	0.016	–	–	–	–	–	3.274	9
–	–	–	0.032	–	–	–	–	–	7.670	10
–	–	–	–	0	–	–	–	–	3.253	11
–	–	–	–	2.0	–	–	–	–	3.250	12
–	–	–	–	–	–	256	–	–	3.249	13
–	–	–	–	–	–	512	–	–	3.270	14
–	–	–	–	–	–	1024	–	–	3.321	15
–	–	–	–	–	–	–	500	–	3.274	16
–	–	–	–	–	–	–	1000	–	3.255	17
–	–	–	–	–	–	–	4000	–	3.252	18
–	–	–	–	–	–	–	–	0	3.286	19
–	–	–	–	–	–	–	–	0.2	3.265	20

Table 63: Hyperparameter ablation for NAdamW on 300m on 1x Chinchilla Data

$\beta_1$	$\beta_2$	$\epsilon$	$\eta$	$g_{\text{norm}}$	$\eta_{\text{min}}$	BSZ	warmup	$\lambda$	Loss	Link
0.95	0.98	1e-10	0.008	1	0	128	2000	0.1	3.248	0
0.8	–	–	–	–	–	–	–	–	7.542	1
0.9	–	–	–	–	–	–	–	–	7.054	2
0.98	–	–	–	–	–	–	–	–	3.254	3
–	0.9	–	–	–	–	–	–	–	3.263	4
–	0.95	–	–	–	–	–	–	–	3.256	5
–	–	1e-25	–	–	–	–	–	–	3.250	6
–	–	1e-20	–	–	–	–	–	–	3.250	7
–	–	1e-15	–	–	–	–	–	–	3.250	8
–	–	–	0.016	–	–	–	–	–	7.538	9
–	–	–	0.032	–	–	–	–	–	7.695	10
–	–	–	–	0	–	–	–	–	3.260	11
–	–	–	–	2.0	–	–	–	–	3.250	12
–	–	–	–	–	–	256	–	–	3.272	13
–	–	–	–	–	–	512	–	–	3.344	14
–	–	–	–	–	–	–	500	–	7.744	15
–	–	–	–	–	–	–	1000	–	7.421	16
–	–	–	–	–	–	–	4000	–	3.256	17
–	–	–	–	–	–	–	–	0	3.283	18
–	–	–	–	–	–	–	–	0.2	3.257	19

## C.6 Sweeping Results for Adam-Mini

Table 64: Hyperparameter ablation for Adam-Mini on 130m on 1x Chinchilla Data

$\beta_1$	$\beta_2$	$\epsilon$	$\eta$	$g_{\text{norm}}$	$\eta_{\text{min}}$	BSZ	warmup	$\lambda$	Loss	Link
0.9	0.98	1e-15	0.008	1	0	128	2000	0.1	3.542	0
0.95	–	–	–	–	–	–	–	–	3.560	1
0.98	–	–	–	–	–	–	–	–	7.733	2
–	0.9	–	–	–	–	–	–	–	3.554	3
–	0.95	–	–	–	–	–	–	–	3.545	4
–	–	1e-25	–	–	–	–	–	–	3.546	5
–	–	1e-20	–	–	–	–	–	–	3.546	6
–	–	1e-10	–	–	–	–	–	–	3.548	7
–	–	–	0.004	–	–	–	–	–	3.558	8
–	–	–	0.016	–	–	–	–	–	7.800	9
–	–	–	0.032	–	–	–	–	–	7.825	10
–	–	–	–	0	–	–	–	–	3.542	11
–	–	–	–	2.0	–	–	–	–	3.542	12
–	–	–	–	–	–	256	–	–	3.785	13
–	–	–	–	–	–	–	500	–	7.725	14
–	–	–	–	–	–	–	1000	–	7.729	15
–	–	–	–	–	–	–	4000	–	3.587	16
–	–	–	–	–	–	–	–	0	3.566	17
–	–	–	–	–	–	–	–	0.2	3.589	18

Table 65: Hyperparameter ablation for Adam-Mini on 130m on 2x Chinchilla Data

$\beta_1$	$\beta_2$	$\epsilon$	$\eta$	$g_{\text{norm}}$	$\eta_{\text{min}}$	BSZ	warmup	$\lambda$	Loss	Link
0.9	0.98	1e-20	0.008	2	0	128	2000	0.1	3.416	0
0.95	–	–	–	–	–	–	–	–	3.429	1
0.98	–	–	–	–	–	–	–	–	7.520	2
–	0.9	–	–	–	–	–	–	–	3.422	3
–	0.95	–	–	–	–	–	–	–	3.419	4
–	–	1e-25	–	–	–	–	–	–	3.416	5
–	–	1e-15	–	–	–	–	–	–	3.416	6
–	–	1e-10	–	–	–	–	–	–	3.415	7
–	–	–	0.004	–	–	–	–	–	3.425	8
–	–	–	0.016	–	–	–	–	–	7.796	9
–	–	–	0.032	–	–	–	–	–	7.721	10
–	–	–	–	0	–	–	–	–	3.416	11
–	–	–	–	1.0	–	–	–	–	3.416	12
–	–	–	–	–	–	256	–	–	3.487	13
–	–	–	–	–	–	512	–	–	3.758	14
–	–	–	–	–	–	–	500	–	7.617	15
–	–	–	–	–	–	–	1000	–	7.445	16
–	–	–	–	–	–	–	4000	–	3.424	17
–	–	–	–	–	–	–	–	0	3.447	18
–	–	–	–	–	–	–	–	0.2	3.426	19

Table 66: Hyperparameter ablation for Adam-Mini on 130m on 4x Chinchilla Data

$\beta_1$	$\beta_2$	$\epsilon$	$\eta$	$g_{\text{norm}}$	$\eta_{\text{min}}$	BSZ	warmup	$\lambda$	Loss	Link
0.9	0.98	1e-10	0.008	1	0	128	2000	0.1	3.328	0
0.95	–	–	–	–	–	–	–	–	3.360	1
0.98	–	–	–	–	–	–	–	–	7.771	2
–	0.9	–	–	–	–	–	–	–	3.337	3
–	0.95	–	–	–	–	–	–	–	3.331	4
–	–	1e-25	–	–	–	–	–	–	3.331	5
–	–	1e-20	–	–	–	–	–	–	3.331	6
–	–	1e-15	–	–	–	–	–	–	3.334	7
–	–	–	0.004	–	–	–	–	–	3.334	8
–	–	–	0.016	–	–	–	–	–	7.717	9
–	–	–	0.032	–	–	–	–	–	7.652	10
–	–	–	–	0	–	–	–	–	3.329	11
–	–	–	–	2.0	–	–	–	–	3.329	12
–	–	–	–	–	–	256	–	–	3.363	13
–	–	–	–	–	–	512	–	–	3.447	14
–	–	–	–	–	–	1024	–	–	3.784	15
–	–	–	–	–	–	–	500	–	7.855	16
–	–	–	–	–	–	–	1000	–	7.411	17
–	–	–	–	–	–	–	4000	–	3.331	18
–	–	–	–	–	–	–	–	0	3.364	19
–	–	–	–	–	–	–	–	0.2	3.365	20

Table 67: Hyperparameter ablation for Adam-Mini on 130m on 8x Chinchilla Data

$\beta_1$	$\beta_2$	$\epsilon$	$\eta$	$g_{\text{norm}}$	$\eta_{\text{min}}$	BSZ	warmup	$\lambda$	Loss	Link
0.9	0.98	1e-10	0.008	1	0	128	2000	0.1	3.266	0
0.95	–	–	–	–	–	–	–	–	3.290	1
0.98	–	–	–	–	–	–	–	–	7.552	2
–	0.9	–	–	–	–	–	–	–	3.281	3
–	0.95	–	–	–	–	–	–	–	3.267	4
–	–	1e-25	–	–	–	–	–	–	3.267	5
–	–	1e-20	–	–	–	–	–	–	3.267	6
–	–	1e-15	–	–	–	–	–	–	3.266	7
–	–	–	0.004	–	–	–	–	–	3.264	8
–	–	–	0.016	–	–	–	–	–	7.614	9
–	–	–	0.032	–	–	–	–	–	7.773	10
–	–	–	–	0	–	–	–	–	3.268	11
–	–	–	–	2.0	–	–	–	–	3.270	12
–	–	–	–	–	–	256	–	–	3.281	13
–	–	–	–	–	–	512	–	–	3.324	14
–	–	–	–	–	–	1024	–	–	3.426	15
–	–	–	–	–	–	–	500	–	7.868	16
–	–	–	–	–	–	–	1000	–	7.022	17
–	–	–	–	–	–	–	4000	–	3.266	18
–	–	–	–	–	–	–	–	0	3.304	19
–	–	–	–	–	–	–	–	0.2	3.291	20

Table 68: Hyperparameter ablation for Adam-Mini on 300m on 1x Chinchilla Data

$\beta_1$	$\beta_2$	$\epsilon$	$\eta$	$g_{\text{norm}}$	$\eta_{\text{min}}$	BSZ	warmup	$\lambda$	Loss	Link
0.9	0.98	1e-25	0.004	2	0	128	2000	0.2	3.272	0
0.95	–	–	–	–	–	–	–	–	3.276	1
0.98	–	–	–	–	–	–	–	–	8.024	2
–	0.9	–	–	–	–	–	–	–	3.279	3
–	0.95	–	–	–	–	–	–	–	3.277	4
–	–	1e-25	–	–	–	–	–	–	3.272	5
–	–	1e-20	–	–	–	–	–	–	3.272	6
–	–	1e-15	–	–	–	–	–	–	3.273	7
–	–	1e-10	–	–	–	–	–	–	3.273	8
–	–	–	0.008	–	–	–	–	–	7.859	9
–	–	–	0.016	–	–	–	–	–	7.990	10
–	–	–	0.032	–	–	–	–	–	8.581	11
–	–	–	–	0	–	–	–	–	3.272	12
–	–	–	–	1.0	–	–	–	–	3.272	13
–	–	–	–	–	–	256	–	–	5.691	14
–	–	–	–	–	–	512	–	–	3.629	15
–	–	–	–	–	–	–	500	–	7.020	16
–	–	–	–	–	–	–	1000	–	7.042	17
–	–	–	–	–	–	–	4000	–	3.280	18
–	–	–	–	–	–	–	–	0	3.336	19
–	–	–	–	–	–	–	–	0.1	3.280	20

Table 69: Hyperparameter ablation for Adam-Mini on 520m on 1x Chinchilla Data

$\beta_1$	$\beta_2$	$\epsilon$	$\eta$	$g_{\text{norm}}$	$\eta_{\text{min}}$	BSZ	warmup	$\lambda$	Loss	Link
0.9	0.98	1e-10	0.004	0	0	128	4000	0.1	3.112	0
0.95	–	–	–	–	–	–	–	–	3.113	1
0.98	–	–	–	–	–	–	–	–	7.459	2
–	0.9	–	–	–	–	–	–	–	3.120	3
–	0.95	–	–	–	–	–	–	–	3.115	4
–	–	1e-25	–	–	–	–	–	–	3.115	5
–	–	1e-20	–	–	–	–	–	–	3.115	6
–	–	1e-15	–	–	–	–	–	–	3.112	7
–	–	–	0.008	–	–	–	–	–	7.771	8
–	–	–	0.016	–	–	–	–	–	7.746	9
–	–	–	0.032	–	–	–	–	–	7.778	10
–	–	–	–	1.0	–	–	–	–	3.112	11
–	–	–	–	2.0	–	–	–	–	3.112	12
–	–	–	–	–	–	256	–	–	3.133	13
–	–	–	–	–	–	512	–	–	3.200	14
–	–	–	–	–	–	–	500	–	7.457	15
–	–	–	–	–	–	–	1000	–	7.274	16
–	–	–	–	–	–	–	2000	–	7.289	17
–	–	–	–	–	–	–	–	0	7.792	18
–	–	–	–	–	–	–	–	0.2	3.115	19



## C.7 Sweeping Results for Kron

Table 70: Hyperparameter ablation for Kron on 130m on 1x Chinchilla Data

$\beta_1$	blocksize	$\eta$	$g_{\text{norm}}$	$\eta_{\text{min}}$	NormGrad	Blocking	$Init_{pc}$	$\eta_{pc}$	$p_{pc}$	BSZ	$Step_{pc}$	warmup	$\lambda$	Loss	Link
0.95	256	0.002	1	0	True	True	1	0.2	0.05	128	2000	1000	0.7	3.492	0
0.9	—	—	—	—	—	—	—	—	—	—	—	—	—	3.500	1
0.98	—	—	—	—	—	—	—	—	—	—	—	—	—	3.497	2
—	128	—	—	—	—	—	—	—	—	—	—	—	—	3.492	3
—	512	—	—	—	—	—	—	—	—	—	—	—	—	3.494	4
—	—	0.0005	—	—	—	—	—	—	—	—	—	—	—	3.528	5
—	—	0.001	—	—	—	—	—	—	—	—	—	—	—	3.501	6
—	—	0.004	—	—	—	—	—	—	—	—	—	—	—	3.514	7
—	—	0.008	—	—	—	—	—	—	—	—	—	—	—	7.838	8
—	—	—	0.0	—	—	—	—	—	—	—	—	—	—	3.492	9
—	—	—	2.0	—	—	—	—	—	—	—	—	—	—	3.492	10
—	—	—	—	—	False	—	—	—	—	—	—	—	—	3.493	11
—	—	—	—	—	—	—	—	0.1	—	—	—	—	—	3.498	12
—	—	—	—	—	—	—	—	—	0.1	—	—	—	—	3.493	13
—	—	—	—	—	—	—	—	—	—	256	—	—	—	3.525	14
—	—	—	—	—	—	—	—	—	—	512	—	—	—	3.632	15
—	—	—	—	—	—	—	—	—	—	—	500	—	—	3.503	16
—	—	—	—	—	—	—	—	—	—	—	1000	—	—	3.498	17
—	—	—	—	—	—	—	—	—	—	—	—	2000	—	3.507	18
—	—	—	—	—	—	—	—	—	—	—	—	4000	—	3.583	19
—	—	—	—	—	—	—	—	—	—	—	—	—	0.0	3.519	20
—	—	—	—	—	—	—	—	—	—	—	—	—	0.5	3.491	21
—	—	—	—	—	—	—	—	—	—	—	—	—	0.9	3.495	22

Table 71: Hyperparameter ablation for Kron on 130m on 2x Chinchilla Data

$\beta_1$	blocksize	$\eta$	$g_{\text{norm}}$	$\eta_{\text{min}}$	NormGrad	Blocking	$Init_{pc}$	$\eta_{pc}$	$p_{pc}$	BSZ	$Step_{pc}$	warmup	$\lambda$	Loss	Link
0.95	256	0.002	1	0	True	True	1	0.2	0.05	128	2000	1000	0.5	3.389	0
0.9	—	—	—	—	—	—	—	—	—	—	—	—	—	3.393	1
0.98	—	—	—	—	—	—	—	—	—	—	—	—	—	3.391	2
—	512	—	—	—	—	—	—	—	—	—	—	—	—	3.390	3
—	—	0.0005	—	—	—	—	—	—	—	—	—	—	—	3.409	4
—	—	0.001	—	—	—	—	—	—	—	—	—	—	—	3.391	5
—	—	0.004	—	—	—	—	—	—	—	—	—	—	—	3.410	6
—	—	0.008	—	—	—	—	—	—	—	—	—	—	—	7.458	7
—	—	—	0.0	—	—	—	—	—	—	—	—	—	—	3.389	8
—	—	—	2.0	—	—	—	—	—	—	—	—	—	—	3.389	9
—	—	—	—	—	False	—	—	—	—	—	—	—	—	3.419	10
—	—	—	—	—	—	—	—	0.1	—	—	—	—	—	3.390	11
—	—	—	—	—	—	—	—	—	0.1	—	—	—	—	3.388	12
—	—	—	—	—	—	—	—	—	—	256	—	—	—	3.402	13
—	—	—	—	—	—	—	—	—	—	512	—	—	—	3.438	14
—	—	—	—	—	—	—	—	—	—	1024	—	—	—	3.532	15
—	—	—	—	—	—	—	—	—	—	—	500	—	—	3.393	16
—	—	—	—	—	—	—	—	—	—	—	1000	—	—	3.391	17
—	—	—	—	—	—	—	—	—	—	—	—	2000	—	3.392	18
—	—	—	—	—	—	—	—	—	—	—	—	4000	—	3.407	19
—	—	—	—	—	—	—	—	—	—	—	—	—	0.0	3.420	20
—	—	—	—	—	—	—	—	—	—	—	—	—	0.7	3.392	21
—	—	—	—	—	—	—	—	—	—	—	—	—	0.9	3.401	22

Table 72: Hyperparameter ablation for Kron on 130m on 4x Chinchilla Data

$\beta_1$	blocksize	$\eta$	$g_{\text{norm}}$	$\eta_{\text{min}}$	NormGrad	Blocking	$Init_{pc}$	$\eta_{pc}$	$p_{pc}$	BSZ	$Step_{pc}$	warmup	$\lambda$	Loss	Link
0.95	256	0.002	1	0	True	True	1	0.2	0.05	128	2000	1000	0.5	3.307	0
0.9	–	–	–	–	–	–	–	–	–	–	–	–	–	3.316	1
0.98	–	–	–	–	–	–	–	–	–	–	–	–	–	3.310	2
–	128	–	–	–	–	–	–	–	–	–	–	–	–	3.308	3
–	512	–	–	–	–	–	–	–	–	–	–	–	–	3.313	4
–	–	0.0005	–	–	–	–	–	–	–	–	–	–	–	3.316	5
–	–	0.001	–	–	–	–	–	–	–	–	–	–	–	3.303	6
–	–	0.004	–	–	–	–	–	–	–	–	–	–	–	6.703	7
–	–	0.008	–	–	–	–	–	–	–	–	–	–	–	7.492	8
–	–	–	0.0	–	–	–	–	–	–	–	–	–	–	3.307	9
–	–	–	2.0	–	–	–	–	–	–	–	–	–	–	3.307	10
–	–	–	–	–	False	–	–	–	–	–	–	–	–	5.676	11
–	–	–	–	–	–	–	–	0.1	–	–	–	–	–	3.311	12
–	–	–	–	–	–	–	–	–	0.1	–	–	–	–	3.304	13
–	–	–	–	–	–	–	–	–	–	256	–	–	–	3.307	14
–	–	–	–	–	–	–	–	–	–	512	–	–	–	3.327	15
–	–	–	–	–	–	–	–	–	–	1024	–	–	–	3.370	16
–	–	–	–	–	–	–	–	–	–	–	500	–	–	3.311	17
–	–	–	–	–	–	–	–	–	–	–	1000	–	–	3.310	18
–	–	–	–	–	–	–	–	–	–	–	–	2000	–	3.310	19
–	–	–	–	–	–	–	–	–	–	–	–	4000	–	3.317	20
–	–	–	–	–	–	–	–	–	–	–	–	–	0.0	3.341	21
–	–	–	–	–	–	–	–	–	–	–	–	–	0.7	3.326	22
–	–	–	–	–	–	–	–	–	–	–	–	–	0.9	3.340	23

Table 73: Hyperparameter ablation for Kron on 130m on 8x Chinchilla Data

$\beta_1$	blocksize	$\eta$	$g_{\text{norm}}$	$\eta_{\text{min}}$	NormGrad	Blocking	$Init_{pc}$	$\eta_{pc}$	$p_{pc}$	BSZ	$Step_{pc}$	warmup	$\lambda$	Loss	Link
0.95	256	0.001	1	0	True	True	1	0.2	0.1	128	2000	1000	0.5	3.239	0
0.9	—	—	—	—	—	—	—	—	—	—	—	—	—	3.243	1
0.98	—	—	—	—	—	—	—	—	—	—	—	—	—	3.239	2
—	512	—	—	—	—	—	—	—	—	—	—	—	—	3.239	3
—	—	0.0005	—	—	—	—	—	—	—	—	—	—	—	3.245	4
—	—	0.002	—	—	—	—	—	—	—	—	—	—	—	3.252	5
—	—	0.004	—	—	—	—	—	—	—	—	—	—	—	6.279	6
—	—	0.008	—	—	—	—	—	—	—	—	—	—	—	7.504	7
—	—	—	0.0	—	—	—	—	—	—	—	—	—	—	3.239	8
—	—	—	2.0	—	—	—	—	—	—	—	—	—	—	3.239	9
—	—	—	—	—	False	—	—	—	—	—	—	—	—	3.247	10
—	—	—	—	—	—	False	—	—	—	—	—	—	—	3.240	11
—	—	—	—	—	—	—	—	0.1	—	—	—	—	—	3.242	12
—	—	—	—	—	—	—	—	—	0.05	—	—	—	—	3.240	13
—	—	—	—	—	—	—	—	—	—	256	—	—	—	3.247	14
—	—	—	—	—	—	—	—	—	—	512	—	—	—	3.264	15
—	—	—	—	—	—	—	—	—	—	1024	—	—	—	3.294	16
—	—	—	—	—	—	—	—	—	—	—	500	—	—	3.240	17
—	—	—	—	—	—	—	—	—	—	—	1000	—	—	3.239	18
—	—	—	—	—	—	—	—	—	—	—	—	2000	—	3.239	19
—	—	—	—	—	—	—	—	—	—	—	—	4000	—	3.240	20
—	—	—	—	—	—	—	—	—	—	—	—	—	0.0	3.282	21
—	—	—	—	—	—	—	—	—	—	—	—	—	0.7	3.241	22
—	—	—	—	—	—	—	—	—	—	—	—	—	0.9	3.246	23

Table 74: Hyperparameter ablation for Kron on 300m on 1x Chinchilla Data

$\beta_1$	blocksize	$\eta$	$g_{\text{norm}}$	$\eta_{\text{min}}$	NormGrad	Blocking	$Init_{pc}$	$\eta_{pc}$	$p_{pc}$	BSZ	$Step_{pc}$	warmup	$\lambda$	Loss	Link
0.95	256	0.001	1	0	True	True	1	0.2	0.1	128	2000	1000	0.5	3.244	0
0.9	—	—	—	—	—	—	—	—	—	—	—	—	—	3.248	1
0.98	—	—	—	—	—	—	—	—	—	—	—	—	—	3.244	2
—	512	—	—	—	—	—	—	—	—	—	—	—	—	3.243	3
—	—	0.0005	—	—	—	—	—	—	—	—	—	—	—	3.263	4
—	—	0.002	—	—	—	—	—	—	—	—	—	—	—	3.245	5
—	—	0.004	—	—	—	—	—	—	—	—	—	—	—	7.636	6
—	—	0.008	—	—	—	—	—	—	—	—	—	—	—	6.923	7
—	—	—	0.0	—	—	—	—	—	—	—	—	—	—	3.244	8
—	—	—	2.0	—	—	—	—	—	—	—	—	—	—	3.244	9
—	—	—	—	—	False	—	—	—	—	—	—	—	—	3.253	10
—	—	—	—	—	—	—	—	0.1	—	—	—	—	—	3.247	11
—	—	—	—	—	—	—	—	—	0.05	—	—	—	—	3.245	12
—	—	—	—	—	—	—	—	—	—	256	—	—	—	3.268	13
—	—	—	—	—	—	—	—	—	—	512	—	—	—	3.309	14
—	—	—	—	—	—	—	—	—	—	1024	—	—	—	3.393	15
—	—	—	—	—	—	—	—	—	—	—	500	—	—	3.248	16
—	—	—	—	—	—	—	—	—	—	—	1000	—	—	3.245	17
—	—	—	—	—	—	—	—	—	—	—	—	2000	—	3.249	18
—	—	—	—	—	—	—	—	—	—	—	—	4000	—	3.255	19
—	—	—	—	—	—	—	—	—	—	—	—	—	0.0	3.280	20
—	—	—	—	—	—	—	—	—	—	—	—	—	0.7	3.241	21
—	—	—	—	—	—	—	—	—	—	—	—	—	0.9	3.243	22

Table 75: Hyperparameter ablation for Kron on 520m on 1x Chinchilla Data

$\beta_1$	blocksize	$\eta$	$g_{\text{norm}}$	$\eta_{\text{min}}$	NormGrad	Blocking	$Init_{pc}$	$\eta_{pc}$	$p_{pc}$	BSZ	$Step_{pc}$	warmup	$\lambda$	Loss	Link
0.95	256	0.001	1	0	True	True	1	0.2	0.1	128	2000	1000	0.5	3.084	0
0.9	—	—	—	—	—	—	—	—	—	—	—	—	—	3.088	1
0.98	—	—	—	—	—	—	—	—	—	—	—	—	—	6.663	2
—	512	—	—	—	—	—	—	—	—	—	—	—	—	3.087	3
—	—	0.0005	—	—	—	—	—	—	—	—	—	—	—	3.095	4
—	—	0.002	—	—	—	—	—	—	—	—	—	—	—	6.412	5
—	—	0.004	—	—	—	—	—	—	—	—	—	—	—	6.927	6
—	—	0.008	—	—	—	—	—	—	—	—	—	—	—	7.018	7
—	—	—	0.0	—	—	—	—	—	—	—	—	—	—	3.084	8
—	—	—	2.0	—	—	—	—	—	—	—	—	—	—	3.084	9
—	—	—	—	—	False	—	—	—	—	—	—	—	—	5.487	10
—	—	—	—	—	—	—	—	0.1	—	—	—	—	—	3.088	11
—	—	—	—	—	—	—	—	—	0.05	—	—	—	—	3.093	12
—	—	—	—	—	—	—	—	—	—	256	—	—	—	3.099	13
—	—	—	—	—	—	—	—	—	—	512	—	—	—	3.126	14
—	—	—	—	—	—	—	—	—	—	—	500	—	—	3.087	15
—	—	—	—	—	—	—	—	—	—	—	1000	—	—	3.085	16
—	—	—	—	—	—	—	—	—	—	—	—	2000	—	3.090	17
—	—	—	—	—	—	—	—	—	—	—	—	4000	—	3.094	18
—	—	—	—	—	—	—	—	—	—	—	—	—	0.0	3.127	19
—	—	—	—	—	—	—	—	—	—	—	—	—	0.7	3.088	20
—	—	—	—	—	—	—	—	—	—	—	—	—	0.9	3.089	21

## C.8 Sweeping Results for Soap

Table 76: Hyperparameter ablation for Soap on 130m on 1x Chinchilla Data

$\beta_1$	$\beta_2$	blocksize	$\epsilon$	$\eta$	$g_{\text{norm}}$	$\eta_{\text{min}}$	Blocking	$f_{pc}$	$\beta_{\text{shampoo}}$	BSZ	warmup	$\lambda$	Loss	Link
0.95	0.98	256	1e-15	0.016	1	0	True	1	0.95	128	1000	0.1	3.483	0
0.8	—	—	—	—	—	—	—	—	—	—	—	—	4.868	1
0.9	—	—	—	—	—	—	—	—	—	—	—	—	4.547	2
—	0.9	—	—	—	—	—	—	—	—	—	—	—	3.496	3
—	0.95	—	—	—	—	—	—	—	—	—	—	—	3.491	4
—	0.99	—	—	—	—	—	—	—	—	—	—	—	3.482	5
—	—	128	—	—	—	—	—	—	—	—	—	—	3.489	6
—	—	512	—	—	—	—	—	—	—	—	—	—	3.487	7
—	—	—	1e-20	—	—	—	—	—	—	—	—	—	3.513	8
—	—	—	1e-10	—	—	—	—	—	—	—	—	—	3.483	9
—	—	—	—	0.004	—	—	—	—	—	—	—	—	3.509	10
—	—	—	—	0.008	—	—	—	—	—	—	—	—	3.491	11
—	—	—	—	—	—	—	—	5	—	—	—	—	3.488	12
—	—	—	—	—	—	—	—	—	0.9	—	—	—	3.488	13
—	—	—	—	—	—	—	—	—	0.98	—	—	—	3.489	14
—	—	—	—	—	—	—	—	—	0.99	—	—	—	3.491	15
—	—	—	—	—	—	—	—	—	—	256	—	—	3.523	16
—	—	—	—	—	—	—	—	—	—	512	—	—	3.611	17
—	—	—	—	—	—	—	—	—	—	—	500	—	3.483	18
—	—	—	—	—	—	—	—	—	—	—	2000	—	3.507	19
—	—	—	—	—	—	—	—	—	—	—	—	0	3.508	20
—	—	—	—	—	—	—	—	—	—	—	—	0.2	3.501	21

Table 77: Hyperparameter ablation for Soap on 130m on 2x Chinchilla Data

$\beta_1$	$\beta_2$	blocksize	$\epsilon$	$\eta$	$g_{\text{norm}}$	$\eta_{\text{min}}$	Blocking	$f_{pc}$	$\beta_{\text{shampoo}}$	BSZ	warmup	$\lambda$	Loss	Link
0.95	0.99	256	1e-15	0.016	1	0	False	1	0.98	128	500	0.1	3.376	0
0.8	–	–	–	–	–	–	–	–	–	–	–	–	5.163	1
0.9	–	–	–	–	–	–	–	–	–	–	–	–	3.397	2
–	0.9	–	–	–	–	–	–	–	–	–	–	–	3.393	3
–	0.95	–	–	–	–	–	–	–	–	–	–	–	3.384	4
–	0.98	–	–	–	–	–	–	–	–	–	–	–	3.380	5
–	–	128	–	–	–	–	–	–	–	–	–	–	3.376	6
–	–	512	–	–	–	–	–	–	–	–	–	–	3.376	7
–	–	–	1e-20	–	–	–	–	–	–	–	–	–	3.384	8
–	–	–	1e-10	–	–	–	–	–	–	–	–	–	3.375	9
–	–	–	–	0.004	–	–	–	–	–	–	–	–	3.388	10
–	–	–	–	0.008	–	–	–	–	–	–	–	–	3.374	11
–	–	–	–	–	–	–	–	10	–	–	–	–	3.381	12
–	–	–	–	–	–	–	–	–	0.9	–	–	–	3.378	13
–	–	–	–	–	–	–	–	–	0.95	–	–	–	3.376	14
–	–	–	–	–	–	–	–	–	0.99	–	–	–	3.379	15
–	–	–	–	–	–	–	–	–	–	256	–	–	3.385	16
–	–	–	–	–	–	–	–	–	–	512	–	–	3.414	17
–	–	–	–	–	–	–	–	–	–	1024	–	–	3.479	18
–	–	–	–	–	–	–	–	–	–	–	1000	–	3.379	19
–	–	–	–	–	–	–	–	–	–	–	2000	–	3.385	20
–	–	–	–	–	–	–	–	–	–	–	–	0	3.437	21
–	–	–	–	–	–	–	–	–	–	–	–	0.2	3.395	22



Table 78: Hyperparameter ablation for Soap on 130m on 4x Chinchilla Data

$\beta_1$	$\beta_2$	blocksize	$\epsilon$	$\eta$	$g_{\text{norm}}$	$\eta_{\text{min}}$	Blocking	$f_{pc}$	$\beta_{\text{shampoo}}$	BSZ	warmup	$\lambda$	Loss	Link
0.95	0.99	256	1e-15	0.008	1	0	False	1	0.98	128	500	0.1	3.295	0
0.8	–	–	–	–	–	–	–	–	–	–	–	–	3.328	1
0.9	–	–	–	–	–	–	–	–	–	–	–	–	3.302	2
–	0.9	–	–	–	–	–	–	–	–	–	–	–	3.312	3
–	0.95	–	–	–	–	–	–	–	–	–	–	–	3.303	4
–	0.98	–	–	–	–	–	–	–	–	–	–	–	3.298	5
–	–	128	–	–	–	–	–	–	–	–	–	–	3.295	6
–	–	512	–	–	–	–	–	–	–	–	–	–	3.295	7
–	–	–	1e-20	–	–	–	–	–	–	–	–	–	3.297	8
–	–	–	1e-10	–	–	–	–	–	–	–	–	–	3.295	9
–	–	–	–	0.004	–	–	–	–	–	–	–	–	3.303	10
–	–	–	–	0.016	–	–	–	–	–	–	–	–	3.303	11
–	–	–	–	–	–	–	–	10	–	–	–	–	3.299	12
–	–	–	–	–	–	–	–	–	0.9	–	–	–	3.296	13
–	–	–	–	–	–	–	–	–	0.95	–	–	–	3.294	14
–	–	–	–	–	–	–	–	–	0.99	–	–	–	3.297	15
–	–	–	–	–	–	–	–	–	–	256	–	–	3.305	16
–	–	–	–	–	–	–	–	–	–	512	–	–	3.325	17
–	–	–	–	–	–	–	–	–	–	1024	–	–	3.358	18
–	–	–	–	–	–	–	–	–	–	–	1000	–	3.298	19
–	–	–	–	–	–	–	–	–	–	–	2000	–	3.300	20
–	–	–	–	–	–	–	–	–	–	–	–	0	3.370	21
–	–	–	–	–	–	–	–	–	–	–	–	0.2	3.304	22

Table 79: Hyperparameter ablation for Soap on 130m on 8x Chinchilla Data

$\beta_1$	$\beta_2$	blocksize	$\epsilon$	$\eta$	$g_{\text{norm}}$	$\eta_{\text{min}}$	Blocking	$f_{\text{pc}}$	$\beta_{\text{shampoo}}$	BSZ	warmup	$\lambda$	Loss	Link
0.95	0.99	512	1e-15	0.008	1	0	True	10	0.98	256	1000	0.1	3.239	0
0.8	–	–	–	–	–	–	–	–	–	–	–	–	3.298	1
0.9	–	–	–	–	–	–	–	–	–	–	–	–	3.250	2
–	0.9	–	–	–	–	–	–	–	–	–	–	–	3.251	3
–	0.95	–	–	–	–	–	–	–	–	–	–	–	3.244	4
–	0.98	–	–	–	–	–	–	–	–	–	–	–	3.241	5
–	–	128	–	–	–	–	–	–	–	–	–	–	3.240	6
–	–	256	–	–	–	–	–	–	–	–	–	–	3.242	7
–	–	–	1e-20	–	–	–	–	–	–	–	–	–	3.239	8
–	–	–	1e-10	–	–	–	–	–	–	–	–	–	3.238	9
–	–	–	–	0.004	–	–	–	–	–	–	–	–	3.248	10
–	–	–	–	0.016	–	–	–	–	–	–	–	–	3.249	11
–	–	–	–	–	–	–	–	–	0.9	–	–	–	3.239	12
–	–	–	–	–	–	–	–	–	0.95	–	–	–	3.239	13
–	–	–	–	–	–	–	–	–	0.99	–	–	–	3.241	14
–	–	–	–	–	–	–	–	–	–	128	–	–	3.242	15
–	–	–	–	–	–	–	–	–	–	512	–	–	3.250	16
–	–	–	–	–	–	–	–	–	–	1024	–	–	3.276	17
–	–	–	–	–	–	–	–	–	–	–	500	–	3.240	18
–	–	–	–	–	–	–	–	–	–	–	2000	–	3.240	19
–	–	–	–	–	–	–	–	–	–	–	–	0	3.308	20
–	–	–	–	–	–	–	–	–	–	–	–	0.2	3.243	21

Table 80: Hyperparameter ablation for Soap on 300m on 1x Chinchilla Data

$\beta_1$	$\beta_2$	blocksize	$\epsilon$	$\eta$	$g_{\text{norm}}$	$\eta_{\text{min}}$	Blocking	$f_{\text{pc}}$	$\beta_{\text{shampoo}}$	BSZ	warmup	$\lambda$	Loss	Link
0.95	0.99	512	1e-10	0.008	1	0	True	10	0.9	128	1000	0.1	3.231	0
0.8	–	–	–	–	–	–	–	–	–	–	–	–	4.544	1
0.9	–	–	–	–	–	–	–	–	–	–	–	–	3.251	2
–	0.9	–	–	–	–	–	–	–	–	–	–	–	3.247	3
–	0.95	–	–	–	–	–	–	–	–	–	–	–	3.240	4
–	0.98	–	–	–	–	–	–	–	–	–	–	–	3.234	5
–	–	128	–	–	–	–	–	–	–	–	–	–	3.239	6
–	–	256	–	–	–	–	–	–	–	–	–	–	3.235	7
–	–	–	1e-20	–	–	–	–	–	–	–	–	–	3.233	8
–	–	–	1e-15	–	–	–	–	–	–	–	–	–	3.233	9
–	–	–	–	0.004	–	–	–	–	–	–	–	–	3.239	10
–	–	–	–	0.016	–	–	–	–	–	–	–	–	5.559	11
–	–	–	–	–	–	–	–	–	0.95	–	–	–	3.233	12
–	–	–	–	–	–	–	–	–	0.98	–	–	–	3.235	13
–	–	–	–	–	–	–	–	–	0.99	–	–	–	3.237	14
–	–	–	–	–	–	–	–	–	–	256	–	–	3.248	15
–	–	–	–	–	–	–	–	–	–	512	–	–	3.283	16
–	–	–	–	–	–	–	–	–	–	–	500	–	3.231	17
–	–	–	–	–	–	–	–	–	–	–	2000	–	3.235	18
–	–	–	–	–	–	–	–	–	–	–	–	0	3.268	19
–	–	–	–	–	–	–	–	–	–	–	–	0.2	3.238	20
–	–	–	–	–	–	–	–	–	–	–	–	0.3	3.249	21

Table 81: Hyperparameter ablation for Soap on 520m on 1x Chinchilla Data

$\beta_1$	$\beta_2$	blocksize	$\epsilon$	$\eta$	$g_{\text{norm}}$	$\eta_{\text{min}}$	Blocking	$f_{pc}$	$\beta_{\text{shampoo}}$	BSZ	warmup	$\lambda$	Loss	Link
0.95	0.99	512	1e-10	0.008	1	0	True	10	0.95	128	1000	0.1	3.079	0
0.8	–	–	–	–	–	–	–	–	–	–	–	–	4.630	1
0.9	–	–	–	–	–	–	–	–	–	–	–	–	4.316	2
–	0.9	–	–	–	–	–	–	–	–	–	–	–	3.097	3
–	0.95	–	–	–	–	–	–	–	–	–	–	–	3.090	4
–	0.98	–	–	–	–	–	–	–	–	–	–	–	3.085	5
–	–	128	–	–	–	–	–	–	–	–	–	–	5.395	6
–	–	256	–	–	–	–	–	–	–	–	–	–	4.392	7
–	–	–	1e-20	–	–	–	–	–	–	–	–	–	3.082	8
–	–	–	1e-15	–	–	–	–	–	–	–	–	–	3.081	9
–	–	–	–	0.004	–	–	–	–	–	–	–	–	3.079	10
–	–	–	–	0.016	–	–	–	–	–	–	–	–	5.762	11
–	–	–	–	–	–	–	–	–	0.9	–	–	–	3.080	12
–	–	–	–	–	–	–	–	–	0.98	–	–	–	3.082	13
–	–	–	–	–	–	–	–	–	0.99	–	–	–	3.083	14
–	–	–	–	–	–	–	–	–	–	256	–	–	3.085	15
–	–	–	–	–	–	–	–	–	–	512	–	–	4.215	16
–	–	–	–	–	–	–	–	–	–	–	500	–	4.163	17
–	–	–	–	–	–	–	–	–	–	–	2000	–	3.081	18
–	–	–	–	–	–	–	–	–	–	–	–	0	6.106	19
–	–	–	–	–	–	–	–	–	–	–	–	0.2	3.527	20
–	–	–	–	–	–	–	–	–	–	–	–	0.3	3.118	21

## C.9 Sweeping Results for Muon

Table 82: Hyperparameter ablation for Muon on 130m on 1x Chinchilla Data

$\eta_{adam}$	$\beta_1$	$\beta_2$	Decay(WSD)	$\epsilon$	$\eta$	Schedule	$g_{norm}$	$\eta_{min}$	$\beta_{muon}$	$\epsilon_{muon}$	BSZ	warmup	$\lambda$	Loss	Link
0.0032	0.8	0.98	0.8	1e-15	0.016	linear	1	0	0.95	1e-05	128	0	0.1	3.464	0
0.0016	—	—	—	—	—	—	—	—	—	—	—	—	—	3.474	1
0.0048	—	—	—	—	—	—	—	—	—	—	—	—	—	3.463	2
—	0.9	—	—	—	—	—	—	—	—	—	—	—	—	3.468	3
—	0.95	—	—	—	—	—	—	—	—	—	—	—	—	3.467	4
—	0.98	—	—	—	—	—	—	—	—	—	—	—	—	3.470	5
—	—	0.9	—	—	—	—	—	—	—	—	—	—	—	3.481	6
—	—	0.95	—	—	—	—	—	—	—	—	—	—	—	3.469	7
—	—	—	0.2	—	—	—	—	—	—	—	—	—	—	3.516	8
—	—	—	0.4	—	—	—	—	—	—	—	—	—	—	3.482	9
—	—	—	0.6	—	—	—	—	—	—	—	—	—	—	3.468	10
—	—	—	1.0	—	—	—	—	—	—	—	—	—	—	3.467	11
—	—	—	—	1e-25	—	—	—	—	—	—	—	—	—	3.465	12
—	—	—	—	1e-20	—	—	—	—	—	—	—	—	—	3.465	13
—	—	—	—	1e-10	—	—	—	—	—	—	—	—	—	3.465	14
—	—	—	—	—	0.008	—	—	—	—	—	—	—	—	3.491	15
—	—	—	—	—	0.032	—	—	—	—	—	—	—	—	3.484	16
—	—	—	—	—	0.064	—	—	—	—	—	—	—	—	3.535	17
—	—	—	—	—	0.128	—	—	—	—	—	—	—	—	3.688	18
—	—	—	—	—	—	—	0	—	—	—	—	—	—	3.464	19
—	—	—	—	—	—	—	2.0	—	—	—	—	—	—	3.464	20
—	—	—	—	—	—	—	—	—	0.8	—	—	—	—	3.503	21
—	—	—	—	—	—	—	—	—	0.9	—	—	—	—	3.473	22
—	—	—	—	—	—	—	—	—	0.98	—	—	—	—	3.474	23
—	—	—	—	—	—	—	—	—	—	1e-25	—	—	—	3.465	24
—	—	—	—	—	—	—	—	—	—	1e-20	—	—	—	3.465	25
—	—	—	—	—	—	—	—	—	—	1e-15	—	—	—	3.465	26
—	—	—	—	—	—	—	—	—	—	1e-10	—	—	—	3.465	27
—	—	—	—	—	—	—	—	—	—	—	256	—	—	3.504	28
—	—	—	—	—	—	—	—	—	—	—	512	—	—	3.601	29
—	—	—	—	—	—	—	—	—	—	—	1024	—	—	3.811	30
—	—	—	—	—	—	—	—	—	—	—	—	—	0	3.532	31
—	—	—	—	—	—	—	—	—	—	—	—	—	0.2	3.483	32

Table 83: Hyperparameter ablation for Muon on 130m on 2x Chinchilla Data

$\eta_{adam}$	$\beta_1$	$\beta_2$	Decay(WSD)	$\epsilon$	$\eta$	Schedule	$g_{norm}$	$\eta_{min}$	$\beta_{muon}$	$\epsilon_{muon}$	BSZ	warmup	$\lambda$	Loss	Link
0.0024	0.8	0.98	0.8	1e-15	0.008	linear	1	0	0.98	1e-05	128	0	0.1	3.369	0
0.0008	–	–	–	–	–	–	–	–	–	–	–	–	–	3.385	1
0.0016	–	–	–	–	–	–	–	–	–	–	–	–	–	3.373	2
–	0.9	–	–	–	–	–	–	–	–	–	–	–	–	3.372	3
–	0.95	–	–	–	–	–	–	–	–	–	–	–	–	3.373	4
–	0.98	–	–	–	–	–	–	–	–	–	–	–	–	3.375	5
–	–	0.9	–	–	–	–	–	–	–	–	–	–	–	3.390	6
–	–	0.95	–	–	–	–	–	–	–	–	–	–	–	3.377	7
–	–	–	0.2	–	–	–	–	–	–	–	–	–	–	3.404	8
–	–	–	0.4	–	–	–	–	–	–	–	–	–	–	3.380	9
–	–	–	0.6	–	–	–	–	–	–	–	–	–	–	3.371	10
–	–	–	1.0	–	–	–	–	–	–	–	–	–	–	3.374	11
–	–	–	–	1e-25	–	–	–	–	–	–	–	–	–	3.370	12
–	–	–	–	1e-20	–	–	–	–	–	–	–	–	–	3.370	13
–	–	–	–	1e-10	–	–	–	–	–	–	–	–	–	3.369	14
–	–	–	–	–	0.016	–	–	–	–	–	–	–	–	3.373	15
–	–	–	–	–	0.032	–	–	–	–	–	–	–	–	3.406	16
–	–	–	–	–	0.064	–	–	–	–	–	–	–	–	3.549	17
–	–	–	–	–	0.128	–	–	–	–	–	–	–	–	5.932	18
–	–	–	–	–	–	–	0	–	–	–	–	–	–	3.371	19
–	–	–	–	–	–	–	2.0	–	–	–	–	–	–	3.370	20
–	–	–	–	–	–	–	–	–	0.8	–	–	–	–	3.417	21
–	–	–	–	–	–	–	–	–	0.9	–	–	–	–	3.385	22
–	–	–	–	–	–	–	–	–	0.95	–	–	–	–	3.371	23
–	–	–	–	–	–	–	–	–	–	1e-25	–	–	–	3.370	24
–	–	–	–	–	–	–	–	–	–	1e-20	–	–	–	3.370	25
–	–	–	–	–	–	–	–	–	–	1e-15	–	–	–	3.370	26
–	–	–	–	–	–	–	–	–	–	1e-10	–	–	–	3.370	27
–	–	–	–	–	–	–	–	–	–	–	256	–	–	3.400	28
–	–	–	–	–	–	–	–	–	–	–	512	–	–	3.465	29
–	–	–	–	–	–	–	–	–	–	–	1024	–	–	3.608	30
–	–	–	–	–	–	–	–	–	–	–	–	–	0	3.431	31
–	–	–	–	–	–	–	–	–	–	–	–	–	0.2	3.378	32

Table 84: Hyperparameter ablation for Muon on 130m on 4x Chinchilla Data

$\eta_{adam}$	$\beta_1$	$\beta_2$	Decay(WSD)	$\epsilon$	$\eta$	Schedule	$g_{norm}$	$\eta_{min}$	$\beta_{muon}$	$\epsilon_{muon}$	BSZ	warmup	$\lambda$	Loss	Link
0.0024	0.8	0.98	0.8	1e-15	0.008	linear	1	0	0.95	1e-05	128	0	0.1	3.296	0
0.0008	–	–	–	–	–	–	–	–	–	–	–	–	–	3.309	1
0.0016	–	–	–	–	–	–	–	–	–	–	–	–	–	3.300	2
–	0.9	–	–	–	–	–	–	–	–	–	–	–	–	3.297	3
–	0.95	–	–	–	–	–	–	–	–	–	–	–	–	3.299	4
–	0.98	–	–	–	–	–	–	–	–	–	–	–	–	3.299	5
–	–	0.9	–	–	–	–	–	–	–	–	–	–	–	3.317	6
–	–	0.95	–	–	–	–	–	–	–	–	–	–	–	3.305	7
–	–	–	0.2	–	–	–	–	–	–	–	–	–	–	3.334	8
–	–	–	0.4	–	–	–	–	–	–	–	–	–	–	3.310	9
–	–	–	0.6	–	–	–	–	–	–	–	–	–	–	3.301	10
–	–	–	1.0	–	–	–	–	–	–	–	–	–	–	3.297	11
–	–	–	–	1e-25	–	–	–	–	–	–	–	–	–	3.297	12
–	–	–	–	1e-20	–	–	–	–	–	–	–	–	–	3.296	13
–	–	–	–	1e-10	–	–	–	–	–	–	–	–	–	3.297	14
–	–	–	–	–	0.016	–	–	–	–	–	–	–	–	3.306	15
–	–	–	–	–	0.032	–	–	–	–	–	–	–	–	3.347	16
–	–	–	–	–	0.064	–	–	–	–	–	–	–	–	3.431	17
–	–	–	–	–	0.128	–	–	–	–	–	–	–	–	4.807	18
–	–	–	–	–	–	–	0	–	–	–	–	–	–	3.296	19
–	–	–	–	–	–	–	2.0	–	–	–	–	–	–	3.296	20
–	–	–	–	–	–	–	–	–	0.8	–	–	–	–	3.327	21
–	–	–	–	–	–	–	–	–	0.9	–	–	–	–	3.307	22
–	–	–	–	–	–	–	–	–	0.98	–	–	–	–	3.295	23
–	–	–	–	–	–	–	–	–	–	1e-25	–	–	–	3.296	24
–	–	–	–	–	–	–	–	–	–	1e-20	–	–	–	3.296	25
–	–	–	–	–	–	–	–	–	–	1e-15	–	–	–	3.296	26
–	–	–	–	–	–	–	–	–	–	1e-10	–	–	–	3.297	27
–	–	–	–	–	–	–	–	–	–	–	256	–	–	3.308	28
–	–	–	–	–	–	–	–	–	–	–	512	–	–	3.346	29
–	–	–	–	–	–	–	–	–	–	–	1024	–	–	3.420	30
–	–	–	–	–	–	–	–	–	–	–	–	–	0	3.352	31
–	–	–	–	–	–	–	–	–	–	–	–	–	0.2	3.316	32

Table 85: Hyperparameter ablation for Muon on 130m on 8x Chinchilla Data

$\eta_{adam}$	$\beta_1$	$\beta_2$	Decay(WSD)	$\epsilon$	$\eta$	Schedule	$g_{norm}$	$\eta_{min}$	$\beta_{muon}$	$\epsilon_{muon}$	BSZ	warmup	$\lambda$	Loss	Link
0.0024	0.8	0.98	1	1e-15	0.008	linear	1	0	0.98	1e-05	128	0	0.1	3.240	0
0.0008	–	–	–	–	–	–	–	–	–	–	–	–	–	3.249	1
0.0016	–	–	–	–	–	–	–	–	–	–	–	–	–	3.243	2
–	0.9	–	–	–	–	–	–	–	–	–	–	–	–	3.240	3
–	0.95	–	–	–	–	–	–	–	–	–	–	–	–	3.240	4
–	0.98	–	–	–	–	–	–	–	–	–	–	–	–	3.241	5
–	–	0.9	–	–	–	–	–	–	–	–	–	–	–	3.260	6
–	–	0.95	–	–	–	–	–	–	–	–	–	–	–	3.247	7
–	–	–	0.2	–	–	–	–	–	–	–	–	–	–	3.280	8
–	–	–	0.4	–	–	–	–	–	–	–	–	–	–	3.257	9
–	–	–	0.6	–	–	–	–	–	–	–	–	–	–	3.248	10
–	–	–	0.8	–	–	–	–	–	–	–	–	–	–	3.242	11
–	–	–	–	1e-25	–	–	–	–	–	–	–	–	–	3.238	12
–	–	–	–	1e-20	–	–	–	–	–	–	–	–	–	3.238	13
–	–	–	–	1e-10	–	–	–	–	–	–	–	–	–	3.241	14
–	–	–	–	–	0.016	–	–	–	–	–	–	–	–	3.253	15
–	–	–	–	–	0.032	–	–	–	–	–	–	–	–	3.298	16
–	–	–	–	–	0.064	–	–	–	–	–	–	–	–	3.407	17
–	–	–	–	–	0.128	–	–	–	–	–	–	–	–	5.568	18
–	–	–	–	–	–	–	0	–	–	–	–	–	–	3.240	19
–	–	–	–	–	–	–	2.0	–	–	–	–	–	–	3.239	20
–	–	–	–	–	–	–	–	–	0.8	–	–	–	–	3.265	21
–	–	–	–	–	–	–	–	–	0.9	–	–	–	–	3.250	22
–	–	–	–	–	–	–	–	–	0.95	–	–	–	–	3.243	23
–	–	–	–	–	–	–	–	–	–	1e-25	–	–	–	3.239	24
–	–	–	–	–	–	–	–	–	–	1e-20	–	–	–	3.239	25
–	–	–	–	–	–	–	–	–	–	1e-15	–	–	–	3.239	26
–	–	–	–	–	–	–	–	–	–	1e-10	–	–	–	3.239	27
–	–	–	–	–	–	–	–	–	–	–	256	–	–	3.240	28
–	–	–	–	–	–	–	–	–	–	–	512	–	–	3.257	29
–	–	–	–	–	–	–	–	–	–	–	1024	–	–	3.302	30
–	–	–	–	–	–	–	–	–	–	–	–	–	0	3.310	31
–	–	–	–	–	–	–	–	–	–	–	–	–	0.2	3.256	32

Table 86: Hyperparameter ablation for Muon on 300m on 1x Chinchilla Data

$\eta_{adam}$	$\beta_1$	$\beta_2$	Decay(WSD)	$\epsilon$	$\eta$	Schedule	$g_{norm}$	$\eta_{min}$	$\beta_{muon}$	$\epsilon_{muon}$	BSZ	warmup	$\lambda$	Loss	Link
0.0024	0.8	0.98	0.8	1e-15	0.008	linear	1	0	0.98	1e-05	128	0	0.1	3.224	0
0.0008	–	–	–	–	–	–	–	–	–	–	–	–	–	3.232	1
0.0016	–	–	–	–	–	–	–	–	–	–	–	–	–	3.225	2
–	0.9	–	–	–	–	–	–	–	–	–	–	–	–	3.226	3
–	0.95	–	–	–	–	–	–	–	–	–	–	–	–	3.226	4
–	0.98	–	–	–	–	–	–	–	–	–	–	–	–	3.227	5
–	–	0.9	–	–	–	–	–	–	–	–	–	–	–	3.241	6
–	–	0.95	–	–	–	–	–	–	–	–	–	–	–	3.231	7
–	–	–	0.2	–	–	–	–	–	–	–	–	–	–	3.271	8
–	–	–	0.4	–	–	–	–	–	–	–	–	–	–	3.240	9
–	–	–	0.6	–	–	–	–	–	–	–	–	–	–	3.229	10
–	–	–	1.0	–	–	–	–	–	–	–	–	–	–	3.225	11
–	–	–	–	1e-25	–	–	–	–	–	–	–	–	–	3.224	12
–	–	–	–	1e-20	–	–	–	–	–	–	–	–	–	3.224	13
–	–	–	–	1e-10	–	–	–	–	–	–	–	–	–	3.225	14
–	–	–	–	–	0.016	–	–	–	–	–	–	–	–	3.236	15
–	–	–	–	–	0.032	–	–	–	–	–	–	–	–	3.281	16
–	–	–	–	–	0.064	–	–	–	–	–	–	–	–	4.695	17
–	–	–	–	–	0.128	–	–	–	–	–	–	–	–	5.976	18
–	–	–	–	–	–	–	0	–	–	–	–	–	–	3.224	19
–	–	–	–	–	–	–	2.0	–	–	–	–	–	–	3.225	20
–	–	–	–	–	–	–	–	–	0.8	–	–	–	–	3.257	21
–	–	–	–	–	–	–	–	–	0.9	–	–	–	–	3.235	22
–	–	–	–	–	–	–	–	–	0.95	–	–	–	–	3.224	23
–	–	–	–	–	–	–	–	–	–	1e-25	–	–	–	3.226	24
–	–	–	–	–	–	–	–	–	–	1e-20	–	–	–	3.226	25
–	–	–	–	–	–	–	–	–	–	1e-15	–	–	–	3.226	26
–	–	–	–	–	–	–	–	–	–	1e-10	–	–	–	3.226	27
–	–	–	–	–	–	–	–	–	–	–	256	–	–	3.244	28
–	–	–	–	–	–	–	–	–	–	–	512	–	–	3.301	29
–	–	–	–	–	–	–	–	–	–	–	1024	–	–	3.419	30
–	–	–	–	–	–	–	–	–	–	–	–	–	0	3.285	31
–	–	–	–	–	–	–	–	–	–	–	–	–	0.2	3.234	32



Table 87: Hyperparameter ablation for Muon on 520m on 1x Chinchilla Data

$\eta_{adam}$	$\beta_1$	$\beta_2$	Decay(WSD)	$\epsilon$	$\eta$	Schedule	$g_{norm}$	$\eta_{min}$	$\beta_{muon}$	$\epsilon_{muon}$	BSZ	warmup	$\lambda$	Loss	Link
0.0024	0.8	0.98	1	1e-25	0.008	linear	1	0	0.98	1e-05	128	0	0.1	3.073	0
0.0008	–	–	–	–	–	–	–	–	–	–	–	–	–	3.076	1
0.0016	–	–	–	–	–	–	–	–	–	–	–	–	–	3.072	2
–	0.9	–	–	–	–	–	–	–	–	–	–	–	–	3.073	3
–	0.95	–	–	–	–	–	–	–	–	–	–	–	–	3.074	4
–	0.98	–	–	–	–	–	–	–	–	–	–	–	–	3.073	5
–	–	0.9	–	–	–	–	–	–	–	–	–	–	–	3.090	6
–	–	0.95	–	–	–	–	–	–	–	–	–	–	–	3.080	7
–	–	–	0.2	–	–	–	–	–	–	–	–	–	–	3.134	8
–	–	–	0.4	–	–	–	–	–	–	–	–	–	–	3.101	9
–	–	–	0.6	–	–	–	–	–	–	–	–	–	–	3.086	10
–	–	–	0.8	–	–	–	–	–	–	–	–	–	–	3.076	11
–	–	–	–	1e-20	–	–	–	–	–	–	–	–	–	3.074	12
–	–	–	–	1e-15	–	–	–	–	–	–	–	–	–	3.073	13
–	–	–	–	1e-10	–	–	–	–	–	–	–	–	–	3.073	14
–	–	–	–	–	0.016	–	–	–	–	–	–	–	–	3.095	15
–	–	–	–	–	0.032	–	–	–	–	–	–	–	–	3.147	16
–	–	–	–	–	0.064	–	–	–	–	–	–	–	–	7.886	17
–	–	–	–	–	0.128	–	–	–	–	–	–	–	–	7.900	18
–	–	–	–	–	–	–	0	–	–	–	–	–	–	3.072	19
–	–	–	–	–	–	–	2.0	–	–	–	–	–	–	3.071	20
–	–	–	–	–	–	–	–	–	0.8	–	–	–	–	3.098	21
–	–	–	–	–	–	–	–	–	0.9	–	–	–	–	3.080	22
–	–	–	–	–	–	–	–	–	0.95	–	–	–	–	3.073	23
–	–	–	–	–	–	–	–	–	–	1e-25	–	–	–	3.072	24
–	–	–	–	–	–	–	–	–	–	1e-20	–	–	–	3.072	25
–	–	–	–	–	–	–	–	–	–	1e-15	–	–	–	3.072	26
–	–	–	–	–	–	–	–	–	–	1e-10	–	–	–	3.072	27
–	–	–	–	–	–	–	–	–	–	–	256	–	–	3.080	28
–	–	–	–	–	–	–	–	–	–	–	512	–	–	3.115	29
–	–	–	–	–	–	–	–	–	–	–	1024	–	–	3.188	30
–	–	–	–	–	–	–	–	–	–	–	–	–	0	3.134	31
–	–	–	–	–	–	–	–	–	–	–	–	–	0.2	3.091	32

## C.10 Sweeping Results for Scion

Table 88: Hyperparameter ablation for Scion on 130m on 1x Chinchilla Data

$\eta_{adam}$	$\beta_1$	Decay(WSD)	$\eta$	Schedule	$g_{norm}$	$\eta_{min}$	$\beta_{muon}$	$\epsilon_{scion}$	BSZ	warmup	$\lambda$	Loss	Link
0.0032	0.95	0.8	0.016	linear	1	0	0.95	1e-15	128	0	0.1	3.477	0
0.0016	—	—	—	—	—	—	—	—	—	—	—	3.477	1
0.0048	—	—	—	—	—	—	—	—	—	—	—	3.481	2
0.0064	—	—	—	—	—	—	—	—	—	—	—	3.485	3
0.008	—	—	—	—	—	—	—	—	—	—	—	3.488	4
0.0096	—	—	—	—	—	—	—	—	—	—	—	3.490	5
—	0.8	—	—	—	—	—	—	—	—	—	—	3.493	6
—	0.9	—	—	—	—	—	—	—	—	—	—	3.479	7
—	0.98	—	—	—	—	—	—	—	—	—	—	3.479	8
—	—	0.0	—	—	—	—	—	—	—	—	—	3.779	9
—	—	0.2	—	—	—	—	—	—	—	—	—	3.525	10
—	—	0.4	—	—	—	—	—	—	—	—	—	3.493	11
—	—	0.6	—	—	—	—	—	—	—	—	—	3.480	12
—	—	1.0	—	—	—	—	—	—	—	—	—	3.477	13
—	—	—	0.008	—	—	—	—	—	—	—	—	3.490	14
—	—	—	0.032	—	—	—	—	—	—	—	—	3.509	15
—	—	—	0.064	—	—	—	—	—	—	—	—	5.195	16
—	—	—	0.128	—	—	—	—	—	—	—	—	6.308	17
—	—	—	—	—	0	—	—	—	—	—	—	3.478	18
—	—	—	—	—	2.0	—	—	—	—	—	—	3.478	19
—	—	—	—	—	—	—	0.8	—	—	—	—	3.480	20
—	—	—	—	—	—	—	0.9	—	—	—	—	3.474	21
—	—	—	—	—	—	—	0.98	—	—	—	—	3.496	22
—	—	—	—	—	—	—	—	1e-20	—	—	—	3.477	23
—	—	—	—	—	—	—	—	1e-10	—	—	—	3.477	24
—	—	—	—	—	—	—	—	1e-05	—	—	—	3.477	25
—	—	—	—	—	—	—	—	—	256	—	—	3.505	26
—	—	—	—	—	—	—	—	—	512	—	—	3.583	27
—	—	—	—	—	—	—	—	—	1024	—	—	3.761	28
—	—	—	—	—	—	—	—	—	—	—	0	3.537	29
—	—	—	—	—	—	—	—	—	—	—	0.2	3.495	30

Table 89: Hyperparameter ablation for Scion on 130m on 2x Chinchilla Data

$\eta_{adam}$	$\beta_1$	Decay(WSD)	$\eta$	Schedule	$g_{norm}$	$\eta_{min}$	$\beta_{muon}$	$\epsilon_{scion}$	BSZ	warmup	$\lambda$	Loss	Link
0.0032	0.95	1	0.016	linear	1	0	0.9	1e-15	128	0	0.1	3.379	0
0.0016	–	–	–	–	–	–	–	–	–	–	–	3.378	1
0.0048	–	–	–	–	–	–	–	–	–	–	–	3.382	2
0.0064	–	–	–	–	–	–	–	–	–	–	–	3.384	3
0.008	–	–	–	–	–	–	–	–	–	–	–	3.386	4
0.0096	–	–	–	–	–	–	–	–	–	–	–	3.388	5
–	0.8	–	–	–	–	–	–	–	–	–	–	3.400	6
–	0.9	–	–	–	–	–	–	–	–	–	–	3.383	7
–	0.98	–	–	–	–	–	–	–	–	–	–	3.379	8
–	–	0.0	–	–	–	–	–	–	–	–	–	3.738	9
–	–	0.2	–	–	–	–	–	–	–	–	–	3.437	10
–	–	0.4	–	–	–	–	–	–	–	–	–	3.405	11
–	–	0.6	–	–	–	–	–	–	–	–	–	3.390	12
–	–	0.8	–	–	–	–	–	–	–	–	–	3.383	13
–	–	–	0.008	–	–	–	–	–	–	–	–	3.390	14
–	–	–	0.032	–	–	–	–	–	–	–	–	3.411	15
–	–	–	0.064	–	–	–	–	–	–	–	–	3.513	16
–	–	–	0.128	–	–	–	–	–	–	–	–	6.193	17
–	–	–	–	–	0	–	–	–	–	–	–	3.379	18
–	–	–	–	–	2.0	–	–	–	–	–	–	3.378	19
–	–	–	–	–	–	–	0.8	–	–	–	–	3.385	20
–	–	–	–	–	–	–	0.95	–	–	–	–	3.380	21
–	–	–	–	–	–	–	0.98	–	–	–	–	3.388	22
–	–	–	–	–	–	–	–	1e-20	–	–	–	3.379	23
–	–	–	–	–	–	–	–	1e-10	–	–	–	3.380	24
–	–	–	–	–	–	–	–	1e-05	–	–	–	3.380	25
–	–	–	–	–	–	–	–	–	256	–	–	3.391	26
–	–	–	–	–	–	–	–	–	512	–	–	3.436	27
–	–	–	–	–	–	–	–	–	1024	–	–	3.532	28
–	–	–	–	–	–	–	–	–	–	–	0	3.443	29
–	–	–	–	–	–	–	–	–	–	–	0.2	3.401	30

Table 90: Hyperparameter ablation for Scion on 130m on 4x Chinchilla Data

$\eta_{adam}$	$\beta_1$	Decay(WSD)	$\eta$	Schedule	$g_{norm}$	$\eta_{min}$	$\beta_{muon}$	$\epsilon_{scion}$	BSZ	warmup	$\lambda$	Loss	Link
0.0016	0.98	1	0.008	linear	2	0	0.9	1e-15	128	0	0.1	3.302	0
0.0008	–	–	–	–	–	–	–	–	–	–	–	3.302	1
0.0024	–	–	–	–	–	–	–	–	–	–	–	3.306	2
0.0032	–	–	–	–	–	–	–	–	–	–	–	3.307	3
0.004	–	–	–	–	–	–	–	–	–	–	–	3.306	4
0.0048	–	–	–	–	–	–	–	–	–	–	–	3.307	5
–	0.8	–	–	–	–	–	–	–	–	–	–	3.331	6
–	0.9	–	–	–	–	–	–	–	–	–	–	3.314	7
–	0.95	–	–	–	–	–	–	–	–	–	–	3.306	8
–	–	0.0	–	–	–	–	–	–	–	–	–	3.584	9
–	–	0.2	–	–	–	–	–	–	–	–	–	3.340	10
–	–	0.4	–	–	–	–	–	–	–	–	–	3.317	11
–	–	0.6	–	–	–	–	–	–	–	–	–	3.307	12
–	–	0.8	–	–	–	–	–	–	–	–	–	3.303	13
–	–	–	0.016	–	–	–	–	–	–	–	–	3.310	14
–	–	–	0.032	–	–	–	–	–	–	–	–	3.361	15
–	–	–	0.064	–	–	–	–	–	–	–	–	5.152	16
–	–	–	0.128	–	–	–	–	–	–	–	–	5.894	17
–	–	–	–	–	0	–	–	–	–	–	–	3.303	18
–	–	–	–	–	1.0	–	–	–	–	–	–	3.303	19
–	–	–	–	–	–	–	0.8	–	–	–	–	3.307	20
–	–	–	–	–	–	–	0.95	–	–	–	–	3.302	21
–	–	–	–	–	–	–	0.98	–	–	–	–	3.303	22
–	–	–	–	–	–	–	–	1e-20	–	–	–	3.302	23
–	–	–	–	–	–	–	–	1e-10	–	–	–	3.303	24
–	–	–	–	–	–	–	–	1e-05	–	–	–	3.302	25
–	–	–	–	–	–	–	–	–	256	–	–	3.323	26
–	–	–	–	–	–	–	–	–	512	–	–	3.365	27
–	–	–	–	–	–	–	–	–	1024	–	–	3.447	28
–	–	–	–	–	–	–	–	–	–	–	0	3.375	29
–	–	–	–	–	–	–	–	–	–	–	0.2	3.309	30

Table 91: Hyperparameter ablation for Scion on 130m on 8x Chinchilla Data

$\eta_{adam}$	$\beta_1$	Decay(WSD)	$\eta$	Schedule	$g_{norm}$	$\eta_{min}$	$\beta_{muon}$	$\epsilon_{scion}$	BSZ	warmup	$\lambda$	Loss	Link
0.0016	0.98	1	0.008	linear	2	0	0.9	1e-05	128	0	0.1	3.246	0
0.0008	–	–	–	–	–	–	–	–	–	–	–	3.245	1
0.0024	–	–	–	–	–	–	–	–	–	–	–	3.250	2
0.0032	–	–	–	–	–	–	–	–	–	–	–	3.249	3
0.004	–	–	–	–	–	–	–	–	–	–	–	3.251	4
0.0048	–	–	–	–	–	–	–	–	–	–	–	3.252	5
–	0.8	–	–	–	–	–	–	–	–	–	–	3.275	6
–	0.9	–	–	–	–	–	–	–	–	–	–	3.257	7
–	0.95	–	–	–	–	–	–	–	–	–	–	3.249	8
–	–	0.0	–	–	–	–	–	–	–	–	–	3.558	9
–	–	0.2	–	–	–	–	–	–	–	–	–	3.288	10
–	–	0.4	–	–	–	–	–	–	–	–	–	3.266	11
–	–	0.6	–	–	–	–	–	–	–	–	–	3.255	12
–	–	0.8	–	–	–	–	–	–	–	–	–	3.249	13
–	–	–	0.016	–	–	–	–	–	–	–	–	3.265	14
–	–	–	0.032	–	–	–	–	–	–	–	–	3.325	15
–	–	–	0.064	–	–	–	–	–	–	–	–	5.264	16
–	–	–	0.128	–	–	–	–	–	–	–	–	6.174	17
–	–	–	–	–	0	–	–	–	–	–	–	3.247	18
–	–	–	–	–	1.0	–	–	–	–	–	–	3.247	19
–	–	–	–	–	–	–	0.8	–	–	–	–	3.249	20
–	–	–	–	–	–	–	0.95	–	–	–	–	3.245	21
–	–	–	–	–	–	–	0.98	–	–	–	–	3.247	22
–	–	–	–	–	–	–	–	1e-20	–	–	–	3.246	23
–	–	–	–	–	–	–	–	1e-15	–	–	–	3.246	24
–	–	–	–	–	–	–	–	1e-10	–	–	–	3.247	25
–	–	–	–	–	–	–	–	–	256	–	–	3.250	26
–	–	–	–	–	–	–	–	–	512	–	–	3.272	27
–	–	–	–	–	–	–	–	–	1024	–	–	3.319	28
–	–	–	–	–	–	–	–	–	–	–	0	3.308	29
–	–	–	–	–	–	–	–	–	–	–	0.2	3.261	30

Table 92: Hyperparameter ablation for Scion on 300m on 1x Chinchilla Data

$\eta_{adam}$	$\beta_1$	Decay(WSD)	$\eta$	Schedule	$g_{norm}$	$\eta_{min}$	$\beta_{muon}$	$\epsilon_{scion}$	BSZ	warmup	$\lambda$	Loss	Link
0.0008	0.98	0.8	0.008	linear	2	0	0.9	1e-05	128	0	0.1	3.232	0
0.0016	–	–	–	–	–	–	–	–	–	–	–	3.237	1
0.0024	–	–	–	–	–	–	–	–	–	–	–	3.242	2
0.0032	–	–	–	–	–	–	–	–	–	–	–	3.242	3
0.004	–	–	–	–	–	–	–	–	–	–	–	3.241	4
0.0048	–	–	–	–	–	–	–	–	–	–	–	3.243	5
–	0.8	–	–	–	–	–	–	–	–	–	–	3.256	6
–	0.9	–	–	–	–	–	–	–	–	–	–	3.240	7
–	0.95	–	–	–	–	–	–	–	–	–	–	3.236	8
–	–	0.0	–	–	–	–	–	–	–	–	–	3.497	9
–	–	0.2	–	–	–	–	–	–	–	–	–	3.271	10
–	–	0.4	–	–	–	–	–	–	–	–	–	3.245	11
–	–	0.6	–	–	–	–	–	–	–	–	–	3.234	12
–	–	1.0	–	–	–	–	–	–	–	–	–	3.237	13
–	–	–	0.016	–	–	–	–	–	–	–	–	3.242	14
–	–	–	0.032	–	–	–	–	–	–	–	–	3.301	15
–	–	–	0.064	–	–	–	–	–	–	–	–	3.400	16
–	–	–	0.128	–	–	–	–	–	–	–	–	5.627	17
–	–	–	–	–	0	–	–	–	–	–	–	3.233	18
–	–	–	–	–	1.0	–	–	–	–	–	–	3.233	19
–	–	–	–	–	–	–	0.8	–	–	–	–	3.240	20
–	–	–	–	–	–	–	0.95	–	–	–	–	3.231	21
–	–	–	–	–	–	–	0.98	–	–	–	–	3.236	22
–	–	–	–	–	–	–	–	1e-20	–	–	–	3.232	23
–	–	–	–	–	–	–	–	1e-15	–	–	–	3.232	24
–	–	–	–	–	–	–	–	1e-10	–	–	–	3.233	25
–	–	–	–	–	–	–	–	–	256	–	–	3.259	26
–	–	–	–	–	–	–	–	–	512	–	–	3.319	27
–	–	–	–	–	–	–	–	–	1024	–	–	3.432	28
–	–	–	–	–	–	–	–	–	–	–	0	3.303	29
–	–	–	–	–	–	–	–	–	–	–	0.2	3.241	30

Table 93: Hyperparameter ablation for Scion on 520m on 1x Chinchilla Data

$\eta_{adam}$	$\beta_1$	Decay(WSD)	$\eta$	Schedule	$g_{norm}$	$\eta_{min}$	$\beta_{muon}$	$\epsilon_{scion}$	BSZ	warmup	$\lambda$	Loss	Link
0.0008	0.98	1	0.008	linear	2	0	0.9	1e-05	128	0	0.1	3.080	0
0.0016	–	–	–	–	–	–	–	–	–	–	–	3.089	1
0.0024	–	–	–	–	–	–	–	–	–	–	–	3.090	2
0.0032	–	–	–	–	–	–	–	–	–	–	–	3.090	3
0.004	–	–	–	–	–	–	–	–	–	–	–	3.090	4
0.0048	–	–	–	–	–	–	–	–	–	–	–	3.091	5
–	0.8	–	–	–	–	–	–	–	–	–	–	3.104	6
–	0.9	–	–	–	–	–	–	–	–	–	–	3.090	7
–	0.95	–	–	–	–	–	–	–	–	–	–	3.081	8
–	–	0.0	–	–	–	–	–	–	–	–	–	3.400	9
–	–	0.2	–	–	–	–	–	–	–	–	–	3.136	10
–	–	0.4	–	–	–	–	–	–	–	–	–	3.106	11
–	–	0.6	–	–	–	–	–	–	–	–	–	3.091	12
–	–	0.8	–	–	–	–	–	–	–	–	–	3.083	13
–	–	–	0.016	–	–	–	–	–	–	–	–	3.105	14
–	–	–	0.032	–	–	–	–	–	–	–	–	3.173	15
–	–	–	0.064	–	–	–	–	–	–	–	–	4.420	16
–	–	–	0.128	–	–	–	–	–	–	–	–	7.335	17
–	–	–	–	–	0	–	–	–	–	–	–	3.082	18
–	–	–	–	–	1.0	–	–	–	–	–	–	3.081	19
–	–	–	–	–	–	–	0.8	–	–	–	–	3.084	20
–	–	–	–	–	–	–	0.95	–	–	–	–	3.079	21
–	–	–	–	–	–	–	0.98	–	–	–	–	3.086	22
–	–	–	–	–	–	–	–	1e-20	–	–	–	3.081	23
–	–	–	–	–	–	–	–	1e-15	–	–	–	3.081	24
–	–	–	–	–	–	–	–	1e-10	–	–	–	3.080	25
–	–	–	–	–	–	–	–	–	256	–	–	3.093	26
–	–	–	–	–	–	–	–	–	512	–	–	3.134	27
–	–	–	–	–	–	–	–	–	1024	–	–	3.215	28
–	–	–	–	–	–	–	–	–	–	–	0	3.150	29
–	–	–	–	–	–	–	–	–	–	–	0.2	3.097	30

## C.11 Sweeping Results for Sophia

Table 94: Hyperparameter ablation for Sophia on 130m on 1x Chinchilla Data

$\beta_1$	$\beta_2$	$\epsilon$	$\gamma$	$\eta$	BSZ	warmup	$\lambda$	Loss	Link
0.95	0.9	1e-07	0.0125	0.004	128	4000	0	3.544	0
0.8	—	—	—	—	—	—	—	3.581	1
0.9	—	—	—	—	—	—	—	3.546	2
0.98	—	—	—	—	—	—	—	3.624	3
—	0.95	—	—	—	—	—	—	3.546	1
—	0.98	—	—	—	—	—	—	3.553	2
—	0.99	—	—	—	—	—	—	3.563	3
—	0.995	—	—	—	—	—	—	3.573	4
—	—	1e-17	—	—	—	—	—	3.559	1
—	—	1e-12	—	—	—	—	—	3.559	2
—	—	—	0.00625	—	—	—	—	3.546	1
—	—	—	0.025	—	—	—	—	3.546	2
—	—	—	0.05	—	—	—	—	3.554	3
—	—	—	—	0.002	—	—	—	3.551	1
—	—	—	—	0.008	—	—	—	3.576	2
—	—	—	—	0.016	—	—	—	7.769	3
—	—	—	—	0.032	—	—	—	7.745	4
—	—	—	—	—	—	500	—	7.754	1
—	—	—	—	—	—	1000	—	7.824	2
—	—	—	—	—	—	2000	—	3.576	3
—	—	—	—	—	—	—	0.1	3.553	1
—	—	—	—	—	—	—	0.2	3.566	2
—	—	—	—	—	—	—	0.3	3.579	3



Table 95: Hyperparameter ablation for Sophia on 130m on 2x Chinchilla Data

$\beta_1$	$\beta_2$	$\epsilon$	$\gamma$	$\eta$	BSZ	warmup	$\lambda$	Loss	Link
0.95	0.9	1e-07	0.0125	0.004	128	4000	0.1	3.414	0
0.8	–	–	–	–	–	–	–	3.439	1
0.9	–	–	–	–	–	–	–	3.416	2
0.98	–	–	–	–	–	–	–	3.469	3
–	0.95	–	–	–	–	–	–	3.414	1
–	0.98	–	–	–	–	–	–	3.418	2
–	0.99	–	–	–	–	–	–	3.421	3
–	0.995	–	–	–	–	–	–	3.428	4
–	–	1e-17	–	–	–	–	–	3.426	1
–	–	1e-12	–	–	–	–	–	3.424	2
–	–	–	0.00625	–	–	–	–	3.416	1
–	–	–	0.025	–	–	–	–	3.415	2
–	–	–	0.05	–	–	–	–	3.423	3
–	–	–	–	0.002	–	–	–	3.417	1
–	–	–	–	0.008	–	–	–	3.438	2
–	–	–	–	0.016	–	–	–	7.282	3
–	–	–	–	0.032	–	–	–	6.938	4
–	–	–	–	–	256	–	–	3.446	1
–	–	–	–	–	–	500	–	7.811	1
–	–	–	–	–	–	1000	–	7.549	2
–	–	–	–	–	–	2000	–	3.436	3
–	–	–	–	–	–	–	0	3.431	1
–	–	–	–	–	–	–	0.2	3.415	2
–	–	–	–	–	–	–	0.3	3.419	3

Table 96: Hyperparameter ablation for Sophia on 130m on 4x Chinchilla Data

$\beta_1$	$\beta_2$	$\epsilon$	$\gamma$	$\eta$	BSZ	warmup	$\lambda$	Loss	Link
0.95	0.99	1e-07	0.0125	0.004	128	4000	0.2	3.330	0
0.8	–	–	–	–	–	–	–	3.351	1
0.9	–	–	–	–	–	–	–	3.332	2
0.98	–	–	–	–	–	–	–	3.371	3
–	0.9	–	–	–	–	–	–	3.328	1
–	0.95	–	–	–	–	–	–	3.328	2
–	0.98	–	–	–	–	–	–	3.328	3
–	–	1e-17	–	–	–	–	–	3.337	1
–	–	1e-12	–	–	–	–	–	3.338	2
–	–	–	0.00625	–	–	–	–	3.330	1
–	–	–	0.025	–	–	–	–	3.330	2
–	–	–	–	0.002	–	–	–	3.329	1
–	–	–	–	0.016	–	–	–	7.059	3
–	–	–	–	0.032	–	–	–	6.664	4
–	–	–	–	–	256	–	–	3.340	1
–	–	–	–	–	512	–	–	3.390	2
–	–	–	–	–	–	500	–	7.345	1
–	–	–	–	–	–	1000	–	7.022	2
–	–	–	–	–	–	2000	–	3.349	3
–	–	–	–	–	–	–	0	3.367	1
–	–	–	–	–	–	–	0.1	3.330	2
–	–	–	–	–	–	–	0.3	3.332	3

Table 97: Hyperparameter ablation for Sophia on 130m on 8x Chinchilla Data

$\beta_1$	$\beta_2$	$\epsilon$	$\gamma$	$\eta$	BSZ	warmup	$\lambda$	Loss	Link
0.95	0.95	1e-07	0.0125	0.002	128	4000	0.2	3.259	0
0.8	–	–	–	–	–	–	–	3.291	1
0.9	–	–	–	–	–	–	–	3.265	2
–	0.9	–	–	–	–	–	–	3.259	1
–	0.98	–	–	–	–	–	–	3.260	2
–	0.99	–	–	–	–	–	–	3.260	3
–	–	1e-17	–	–	–	–	–	3.266	1
–	–	1e-12	–	–	–	–	–	3.265	2
–	–	–	–	0.004	–	–	–	3.265	1
–	–	–	–	0.008	–	–	–	3.308	2
–	–	–	–	–	256	–	–	3.265	1
–	–	–	–	–	–	1000	–	3.277	1
–	–	–	–	–	–	2000	–	3.260	2
–	–	–	–	–	–	–	0	3.300	1

Table 98: Hyperparameter ablation for Sophia on 300m on 1x Chinchilla Data

$\beta_1$	$\beta_2$	$\epsilon$	$\gamma$	$\eta$	BSZ	warmup	$\lambda$	Loss	Link
0.9	0.9	1e-07	0.0125	0.004	128	4000	0.1	3.267	0
0.8	–	–	–	–	–	–	–	3.288	1
0.95	–	–	–	–	–	–	–	7.390	2
0.98	–	–	–	–	–	–	–	7.525	3
–	0.95	–	–	–	–	–	–	3.270	1
–	0.98	–	–	–	–	–	–	3.275	2
–	0.99	–	–	–	–	–	–	3.280	3
–	0.995	–	–	–	–	–	–	3.280	4
–	–	1e-17	–	–	–	–	–	3.289	1
–	–	1e-12	–	–	–	–	–	3.289	2
–	–	–	0.00625	–	–	–	–	3.273	1
–	–	–	0.025	–	–	–	–	3.271	2
–	–	–	0.05	–	–	–	–	3.274	3
–	–	–	–	0.002	–	–	–	3.276	1
–	–	–	–	0.008	–	–	–	6.966	2
–	–	–	–	0.016	–	–	–	7.142	3
–	–	–	–	0.032	–	–	–	6.811	4
–	–	–	–	–	256	–	–	3.298	1
–	–	–	–	–	–	500	–	7.149	1
–	–	–	–	–	–	1000	–	7.093	2
–	–	–	–	–	–	2000	–	7.382	3
–	–	–	–	–	–	–	0	3.288	1
–	–	–	–	–	–	–	0.2	3.274	2
–	–	–	–	–	–	–	0.3	3.281	3

Table 99: Hyperparameter ablation for Sophia on 520m on 1x Chinchilla Data

$\beta_1$	$\beta_2$	$\epsilon$	$\gamma$	$\eta$	BSZ	warmup	$\lambda$	Loss	Link
0.95	0.9	1e-07	0.0125	0.002	128	4000	0.3	3.106	0
0.8	–	–	–	–	–	–	–	3.125	1
0.9	–	–	–	–	–	–	–	3.109	2
0.98	–	–	–	–	–	–	–	3.133	3
–	0.95	–	–	–	–	–	–	3.111	1
–	0.98	–	–	–	–	–	–	3.111	2
–	0.99	–	–	–	–	–	–	6.571	3
–	–	1e-17	–	–	–	–	–	3.116	1
–	–	1e-12	–	–	–	–	–	3.116	2
–	–	–	0.00625	–	–	–	–	3.107	1
–	–	–	0.025	–	–	–	–	3.107	2
–	–	–	–	0.004	–	–	–	6.940	1
–	–	–	–	–	–	1000	–	6.908	1
–	–	–	–	–	–	2000	–	6.823	2
–	–	–	–	–	–	–	0	3.148	1
–	–	–	–	–	–	–	0.1	3.113	2
–	–	–	–	–	–	–	0.2	3.105	3

## D Hyperparameter Ablation in Phase II

We reported the results for the optimizers we swept in Phase II. The result is formulated in the same way as in Phase I.

### D.1 Sweeping Results for AdamW

Table 100: Hyperparameter ablation for AdamW on 300m on 2x Chinchilla Data

$\beta_1$	$\beta_2$	$\epsilon$	$\eta$	$g_{\text{norm}}$	$\eta_{\text{min}}$	BSZ	warmup	$\lambda$	Loss	Link
0.9	0.98	1e-15	0.008	1	0	128	2000	0.1	3.166	0
–	–	–	0.004	–	–	–	–	–	3.167	1
–	–	–	–	–	–	256	–	–	3.170	2
–	–	–	–	–	–	–	–	0.2	3.183	3

Table 101: Hyperparameter ablation for AdamW on 300m on 4x Chinchilla Data

$\beta_1$	$\beta_2$	$\epsilon$	$\eta$	$g_{\text{norm}}$	$\eta_{\text{min}}$	BSZ	warmup	$\lambda$	Loss	Link
0.9	0.98	1e-10	0.008	2	0	256	2000	0.1	3.094	0
–	–	–	0.004	–	–	–	–	–	3.101	1
–	–	–	–	–	–	128	–	–	3.103	2
–	–	–	–	–	–	–	–	0.2	3.103	3

Table 102: Hyperparameter ablation for AdamW on 300m on 8x Chinchilla Data

$\beta_1$	$\beta_2$	$\epsilon$	$\eta$	$g_{\text{norm}}$	$\eta_{\text{min}}$	BSZ	warmup	$\lambda$	Loss	Link
0.9	0.98	1e-10	0.008	2	0	256	2000	0.1	3.043	0
–	–	–	0.004	–	–	–	–	–	3.042	1
–	–	–	–	–	–	128	–	–	3.057	2
–	–	–	–	–	–	–	–	0.2	3.059	3

Table 103: Hyperparameter ablation for AdamW on 520m on 2x Chinchilla Data

$\beta_1$	$\beta_2$	$\epsilon$	$\eta$	$g_{\text{norm}}$	$\eta_{\text{min}}$	BSZ	warmup	$\lambda$	Loss	Link
0.9	0.98	1e-10	0.004	1	0	256	1000	0.2	3.023	0
–	–	–	–	–	–	128	–	–	6.654	1
–	–	–	–	–	–	–	–	0.1	3.025	2

Table 104: Hyperparameter ablation for AdamW on 520m on 4x Chinchilla Data

$\beta_1$	$\beta_2$	$\epsilon$	$\eta$	$g_{\text{norm}}$	$\eta_{\text{min}}$	BSZ	warmup	$\lambda$	Loss	Link
0.9	0.98	1e-10	0.004	1	0	256	1000	0.1	2.958	0
–	–	–	0.008	–	–	–	–	–	7.075	1
–	–	–	–	–	–	128	–	–	7.139	2
–	–	–	–	–	–	–	–	0.2	2.962	3

Table 105: Hyperparameter ablation for AdamW on 520m on 8x Chinchilla Data

$\beta_1$	$\beta_2$	$\epsilon$	$\eta$	$g_{\text{norm}}$	$\eta_{\text{min}}$	BSZ	warmup	$\lambda$	Loss	Link
0.9	0.98	1e-10	0.004	1	0	256	1000	0.1	2.913	0
–	–	–	0.008	–	–	–	–	–	7.183	1
–	–	–	–	–	–	128	–	–	6.932	2
–	–	–	–	–	–	–	–	0.2	2.921	3

## D.2 Sweeping Results for Cautious

Table 106: Hyperparameter ablation for Cautious on 300m on 2x Chinchilla Data

$\beta_1$	$\beta_2$	$\epsilon$	$\eta$	$g_{\text{norm}}$	$\eta_{\text{min}}$	BSZ	warmup	$\lambda$	Loss	Link
0.98	0.98	1e-25	0.008	1	0	256	2000	0.1	3.165	0
–	–	–	0.004	–	–	–	–	–	3.175	1
–	–	–	–	–	–	128	–	–	3.171	2

Table 107: Hyperparameter ablation for Cautious on 300m on 4x Chinchilla Data

$\beta_1$	$\beta_2$	$\epsilon$	$\eta$	$g_{\text{norm}}$	$\eta_{\text{min}}$	BSZ	warmup	$\lambda$	Loss	Link
0.98	0.98	1e-25	0.008	1	0	256	2000	0.1	3.094	0
–	–	–	0.004	–	–	–	–	–	3.098	1
–	–	–	–	–	–	128	–	–	3.114	2

Table 108: Hyperparameter ablation for Cautious on 300m on 8x Chinchilla Data

$\beta_1$	$\beta_2$	$\epsilon$	$\eta$	$g_{\text{norm}}$	$\eta_{\text{min}}$	BSZ	warmup	$\lambda$	Loss	Link
0.98	0.98	1e-25	0.008	1	0	256	2000	0.1	3.043	0
–	–	–	0.004	–	–	–	–	–	3.041	1
–	–	–	–	–	–	128	–	–	3.071	2

Table 109: Hyperparameter ablation for Cautious on 520m on 2x Chinchilla Data

$\beta_1$	$\beta_2$	$\epsilon$	$\eta$	$g_{\text{norm}}$	$\eta_{\text{min}}$	BSZ	warmup	$\lambda$	Loss	Link
0.98	0.98	1e-25	0.008	1	0	256	2000	0.1	3.017	0
–	–	–	0.004	–	–	–	–	–	3.019	1
–	–	–	–	–	–	128	–	–	>10	2

Table 110: Hyperparameter ablation for Cautious on 520m on 4x Chinchilla Data

$\beta_1$	$\beta_2$	$\epsilon$	$\eta$	$g_{\text{norm}}$	$\eta_{\text{min}}$	BSZ	warmup	$\lambda$	Loss	Link
0.98	0.98	1e-25	0.004	1	0	256	2000	0.1	2.956	0
–	–	–	0.008	–	–	–	–	–	2.959	1
–	–	–	–	–	–	128	–	–	2.971	2

Table 111: Hyperparameter ablation for Cautious on 520m on 8x Chinchilla Data

$\beta_1$	$\beta_2$	$\epsilon$	$\eta$	$g_{\text{norm}}$	$\eta_{\text{min}}$	BSZ	warmup	$\lambda$	Loss	Link
0.98	0.98	1e-25	0.004	1	0	256	2000	0.1	2.910	0
–	–	–	0.008	–	–	–	–	–	2.919	1
–	–	–	–	–	–	128	–	–	2.931	2

### D.3 Sweeping Results for Lion

Table 112: Hyperparameter ablation for Lion on 300m on 2x Chinchilla Data

$\beta_1$	$\beta_2$	$\eta$	$g_{\text{norm}}$	$\eta_{\text{min}}$	BSZ	warmup	$\lambda$	Loss	Link
0.9	0.98	0.001	1	0	128	2000	0.7	3.170	0
–	0.9	–	–	–	–	–	–	3.189	1
–	0.95	–	–	–	–	–	–	3.175	2
–	–	0.0005	–	–	–	–	–	3.172	3
–	–	–	–	–	256	–	–	3.183	4

Table 113: Hyperparameter ablation for Lion on 300m on 4x Chinchilla Data

$\beta_1$	$\beta_2$	$\eta$	$g_{\text{norm}}$	$\eta_{\text{min}}$	BSZ	warmup	$\lambda$	Loss	Link
0.9	0.98	0.001	1	0	256	2000	0.7	3.100	0
–	0.9	–	–	–	–	–	–	3.114	1
–	0.95	–	–	–	–	–	–	3.103	2
–	–	0.0005	–	–	–	–	–	3.105	3
–	–	–	–	–	128	–	–	3.104	4

Table 114: Hyperparameter ablation for Lion on 300m on 8x Chinchilla Data

$\beta_1$	$\beta_2$	$\eta$	$g_{\text{norm}}$	$\eta_{\text{min}}$	BSZ	warmup	$\lambda$	Loss	Link
0.9	0.98	0.001	1	0	256	2000	0.7	3.046	0
–	0.9	–	–	–	–	–	–	3.058	1
–	0.95	–	–	–	–	–	–	3.050	2
–	–	0.0005	–	–	–	–	–	3.043	3
–	–	–	–	–	128	–	–	3.061	4

Table 115: Hyperparameter ablation for Lion on 520m on 2x Chinchilla Data

$\beta_1$	$\beta_2$	$\eta$	$g_{\text{norm}}$	$\eta_{\text{min}}$	BSZ	warmup	$\lambda$	Loss	Link
0.9	0.95	0.001	1	0	128	2000	0.6	3.029	0
–	0.9	–	–	–	–	–	–	3.045	1
–	0.98	–	–	–	–	–	–	7.779	2
–	–	0.0005	–	–	–	–	–	3.028	3
–	–	–	–	–	256	–	–	3.030	4

Table 116: Hyperparameter ablation for Lion on 520m on 4x Chinchilla Data

$\beta_1$	$\beta_2$	$\eta$	$g_{\text{norm}}$	$\eta_{\text{min}}$	BSZ	warmup	$\lambda$	Loss	Link
0.9	0.95	0.001	1	0	256	2000	0.6	2.965	0
–	0.9	–	–	–	–	–	–	2.971	1
–	0.98	–	–	–	–	–	–	2.965	2
–	–	0.0005	–	–	–	–	–	2.966	3
–	–	–	–	–	128	–	–	2.975	4

Table 117: Hyperparameter ablation for Lion on 520m on 8x Chinchilla Data

$\beta_1$	$\beta_2$	$\eta$	$g_{\text{norm}}$	$\eta_{\text{min}}$	BSZ	warmup	$\lambda$	Loss	Link
0.9	0.95	0.0005	1	0	256	2000	0.6	2.915	0
–	0.9	–	–	–	–	–	–	2.922	1
–	0.98	–	–	–	–	–	–	2.915	2
–	–	0.001	–	–	–	–	–	2.920	3
–	–	–	–	–	128	–	–	2.922	4

#### D.4 Sweeping Results for Mars

Table 118: Hyperparameter ablation for Mars on 300m on 2x Chinchilla Data

$\beta_1$	$\beta_2$	$\epsilon$	$\gamma$	$\eta$	$g_{\text{norm}}$	$\eta_{\text{min}}$	BSZ	warmup	$\lambda$	Loss	Link
0.98	0.99	1e-25	0.05	0.008	1	0	128	1000	0.1	3.158	0
0.9	–	–	–	–	–	–	–	–	–	3.174	1
0.95	–	–	–	–	–	–	–	–	–	3.156	2
–	–	–	–	0.004	–	–	–	–	–	3.168	3
–	–	–	–	–	–	–	256	–	–	3.171	4
–	–	–	–	–	–	–	–	2000	–	3.159	5
–	–	–	–	–	–	–	–	4000	–	3.167	6

Table 119: Hyperparameter ablation for Mars on 300m on 4x Chinchilla Data

$\beta_1$	$\beta_2$	$\epsilon$	$\gamma$	$\eta$	$g_{\text{norm}}$	$\eta_{\text{min}}$	BSZ	warmup	$\lambda$	Loss	Link
0.98	0.99	1e-25	0.05	0.008	1	0	128	1000	0.1	3.097	0
0.9	–	–	–	–	–	–	–	–	–	3.111	1
0.95	–	–	–	–	–	–	–	–	–	3.096	2
–	–	–	–	0.004	–	–	–	–	–	3.095	3
–	–	–	–	–	–	–	256	–	–	3.098	4
–	–	–	–	–	–	–	–	2000	–	3.097	5
–	–	–	–	–	–	–	–	4000	–	3.101	6

Table 120: Hyperparameter ablation for Mars on 300m on 8x Chinchilla Data

$\beta_1$	$\beta_2$	$\epsilon$	$\gamma$	$\eta$	$g_{\text{norm}}$	$\eta_{\text{min}}$	BSZ	warmup	$\lambda$	Loss	Link
0.98	0.99	1e-25	0.05	0.008	1	0	256	1000	0.1	3.040	0
0.9	–	–	–	–	–	–	–	–	–	3.049	1
0.95	–	–	–	–	–	–	–	–	–	3.038	2
–	–	–	–	0.004	–	–	–	–	–	3.046	3
–	–	–	–	–	–	–	128	–	–	3.050	4
–	–	–	–	–	–	–	–	2000	–	3.049	5
–	–	–	–	–	–	–	–	4000	–	3.043	6

Table 121: Hyperparameter ablation for Mars on 520m on 2x Chinchilla Data

$\beta_1$	$\beta_2$	$\epsilon$	$\gamma$	$\eta$	$g_{\text{norm}}$	$\eta_{\text{min}}$	BSZ	warmup	$\lambda$	Loss	Link
0.95	0.98	1e-25	0.025	0.008	1	0	256	2000	0.1	3.015	0
0.9	–	–	–	–	–	–	–	–	–	3.019	1
0.98	–	–	–	–	–	–	–	–	–	3.014	2
–	–	–	–	0.004	–	–	–	–	–	3.019	3
–	–	–	–	–	–	–	128	–	–	3.025	4
–	–	–	–	–	–	–	–	1000	–	3.023	5
–	–	–	–	–	–	–	–	4000	–	3.019	6

Table 122: Hyperparameter ablation for Mars on 520m on 4x Chinchilla Data

$\beta_1$	$\beta_2$	$\epsilon$	$\gamma$	$\eta$	$g_{\text{norm}}$	$\eta_{\text{min}}$	BSZ	warmup	$\lambda$	Loss	Link
0.95	0.98	1e-25	0.025	0.008	1	0	256	2000	0.1	2.955	0
0.9	–	–	–	–	–	–	–	–	–	2.960	1
0.98	–	–	–	–	–	–	–	–	–	2.953	2
–	–	–	–	0.004	–	–	–	–	–	2.953	3
–	–	–	–	–	–	–	128	–	–	2.974	4
–	–	–	–	–	–	–	–	1000	–	2.964	5
–	–	–	–	–	–	–	–	4000	–	2.956	6

Table 123: Hyperparameter ablation for Mars on 520m on 8x Chinchilla Data

$\beta_1$	$\beta_2$	$\epsilon$	$\gamma$	$\eta$	$g_{\text{norm}}$	$\eta_{\text{min}}$	BSZ	warmup	$\lambda$	Loss	Link
0.95	0.98	1e-25	0.025	0.004	1	0	256	2000	0.1	2.906	0
0.9	–	–	–	–	–	–	–	–	–	2.908	1
0.98	–	–	–	–	–	–	–	–	–	2.906	2
–	–	–	–	0.008	–	–	–	–	–	2.917	3
–	–	–	–	–	–	–	128	–	–	2.916	4
–	–	–	–	–	–	–	–	1000	–	2.906	5
–	–	–	–	–	–	–	–	4000	–	2.907	6



## D.5 Sweeping Results for NAdamW

Table 124: Hyperparameter ablation for NAdamW on 300m on 2x Chinchilla Data

$\beta_1$	$\beta_2$	$\epsilon$	$\eta$	$g_{\text{norm}}$	$\eta_{\text{min}}$	BSZ	warmup	$\lambda$	Loss	Link
0.95	0.98	1e-10	0.008	1	0	128	2000	0.1	3.160	0
–	–	–	0.004	–	–	–	–	–	3.160	1
–	–	–	–	–	–	256	–	–	3.165	2

Table 125: Hyperparameter ablation for NAdamW on 300m on 4x Chinchilla Data

$\beta_1$	$\beta_2$	$\epsilon$	$\eta$	$g_{\text{norm}}$	$\eta_{\text{min}}$	BSZ	warmup	$\lambda$	Loss	Link
0.95	0.98	1e-10	0.008	1	0	256	2000	0.1	3.090	0
–	–	–	0.004	–	–	–	–	–	3.097	1
–	–	–	–	–	–	128	–	–	3.098	2

Table 126: Hyperparameter ablation for NAdamW on 300m on 8x Chinchilla Data

$\beta_1$	$\beta_2$	$\epsilon$	$\eta$	$g_{\text{norm}}$	$\eta_{\text{min}}$	BSZ	warmup	$\lambda$	Loss	Link
0.95	0.98	1e-10	0.008	1	0	256	2000	0.1	3.039	0
–	–	–	0.004	–	–	–	–	–	3.039	1
–	–	–	–	–	–	128	–	–	3.055	2

Table 127: Hyperparameter ablation for NAdamW on 520m on 2x Chinchilla Data

$\beta_1$	$\beta_2$	$\epsilon$	$\eta$	$g_{\text{norm}}$	$\eta_{\text{min}}$	BSZ	warmup	$\lambda$	Loss	Link
0.98	0.98	1e-10	0.004	1	0	128	4000	0.1	3.013	0
–	–	–	0.008	–	–	–	–	–	3.023	1
–	–	–	–	–	–	256	–	–	3.020	2

Table 128: Hyperparameter ablation for NAdamW on 520m on 4x Chinchilla Data

$\beta_1$	$\beta_2$	$\epsilon$	$\eta$	$g_{\text{norm}}$	$\eta_{\text{min}}$	BSZ	warmup	$\lambda$	Loss	Link
0.98	0.98	1e-10	0.004	1	0	128	4000	0.1	2.955	0
–	–	–	0.008	–	–	–	–	–	2.971	1
–	–	–	–	–	–	256	–	–	2.954	2

Table 129: Hyperparameter ablation for NAdamW on 520m on 8x Chinchilla Data

$\beta_1$	$\beta_2$	$\epsilon$	$\eta$	$g_{\text{norm}}$	$\eta_{\text{min}}$	BSZ	warmup	$\lambda$	Loss	Link
0.98	0.98	1e-10	0.004	1	0	256	4000	0.1	2.907	0
–	–	–	0.008	–	–	–	–	–	2.910	1
–	–	–	–	–	–	128	–	–	2.913	2

## D.6 Sweeping Results for Adam-Mini

Table 130: Hyperparameter ablation for Adam-Mini on 300m on 2x Chinchilla Data

$\beta_1$	$\beta_2$	$\epsilon$	$\eta$	$g_{\text{norm}}$	$\eta_{\text{min}}$	BSZ	warmup	$\lambda$	Loss	Link
0.9	0.98	1e-25	0.004	2	0	128	2000	0.2	3.178	0
–	–	–	0.002	–	–	–	–	–	3.180	1
–	–	–	–	–	–	256	–	–	6.960	2
–	–	–	–	–	–	–	4000	–	3.183	3
–	–	–	–	–	–	–	–	0.1	3.179	4

Table 131: Hyperparameter ablation for Adam-Mini on 300m on 4x Chinchilla Data

$\beta_1$	$\beta_2$	$\epsilon$	$\eta$	$g_{\text{norm}}$	$\eta_{\text{min}}$	BSZ	warmup	$\lambda$	Loss	Link
0.9	0.98	1e-25	0.004	2	0	128	2000	0.1	3.103	0
–	–	–	0.002	–	–	–	–	–	3.111	1
–	–	–	–	–	–	256	–	–	3.109	2
–	–	–	–	–	–	–	4000	–	3.104	3
–	–	–	–	–	–	–	–	0.2	3.111	4

Table 132: Hyperparameter ablation for Adam-Mini on 300m on 8x Chinchilla Data

$\beta_1$	$\beta_2$	$\epsilon$	$\eta$	$g_{\text{norm}}$	$\eta_{\text{min}}$	BSZ	warmup	$\lambda$	Loss	Link
0.9	0.98	1e-25	0.002	2	0	128	2000	0.2	3.049	0
–	–	–	0.004	–	–	–	–	–	3.064	1
–	–	–	–	–	–	256	–	–	3.052	2
–	–	–	–	–	–	–	4000	–	3.050	3
–	–	–	–	–	–	–	–	0.1	3.051	4

Table 133: Hyperparameter ablation for Adam-Mini on 520m on 2x Chinchilla Data

$\beta_1$	$\beta_2$	$\epsilon$	$\eta$	$g_{\text{norm}}$	$\eta_{\text{min}}$	BSZ	warmup	$\lambda$	Loss	Link
0.9	0.98	1e-10	0.004	1	0	128	4000	0.1	3.027	0
–	–	–	0.002	–	–	–	–	–	3.031	1
–	–	–	–	–	–	256	–	–	3.032	2
–	–	–	–	–	–	–	2000	–	7.359	3
–	–	–	–	–	–	–	–	0.2	3.037	4

Table 134: Hyperparameter ablation for Adam-Mini on 520m on 4x Chinchilla Data

$\beta_1$	$\beta_2$	$\epsilon$	$\eta$	$g_{\text{norm}}$	$\eta_{\text{min}}$	BSZ	warmup	$\lambda$	Loss	Link
0.9	0.98	1e-10	0.004	1	0	128	4000	0.1	2.966	0
–	–	–	0.002	–	–	–	–	–	2.963	1
–	–	–	–	–	–	256	–	–	2.963	2
–	–	–	–	–	–	–	2000	–	7.529	3
–	–	–	–	–	–	–	–	0.2	2.981	4

Table 135: Hyperparameter ablation for Adam-Mini on 520m on 8x Chinchilla Data

$\beta_1$	$\beta_2$	$\epsilon$	$\eta$	$g_{\text{norm}}$	$\eta_{\text{min}}$	BSZ	warmup	$\lambda$	Loss	Link
0.9	0.98	1e-10	0.004	1	0	256	4000	0.1	2.912	0
–	–	–	0.002	–	–	–	–	–	2.918	1
–	–	–	–	–	–	128	–	–	2.921	2
–	–	–	–	–	–	–	2000	–	7.449	3
–	–	–	–	–	–	–	–	0.2	3.025	4

## D.7 Sweeping Results for Muon

Table 136: Hyperparameter ablation for Muon on 300m on 2x Chinchilla Data

$\eta_{\text{adam}}$	$\beta_1$	$\beta_2$	Decay(WSD)	$\epsilon$	$\eta$	Schedule	$g_{\text{norm}}$	$\eta_{\text{min}}$	$\beta_{\text{muon}}$	$\epsilon_{\text{muon}}$	BSZ	warmup	$\lambda$	Loss	Link
0.0024	0.8	0.98	0.8	1e-15	0.008	linear	1	0	0.98	1e-05	128	0	0.1	3.143	0
–	–	–	–	–	0.004	–	–	–	–	–	–	–	–	3.144	1
–	–	–	–	–	–	–	–	–	–	–	256	–	–	3.145	2

Table 137: Hyperparameter ablation for Muon on 300m on 4x Chinchilla Data

$\eta_{\text{adam}}$	$\beta_1$	$\beta_2$	Decay(WSD)	$\epsilon$	$\eta$	Schedule	$g_{\text{norm}}$	$\eta_{\text{min}}$	$\beta_{\text{muon}}$	$\epsilon_{\text{muon}}$	BSZ	warmup	$\lambda$	Loss	Link
0.0012	0.8	0.98	0.8	1e-15	0.004	linear	1	0	0.98	1e-05	128	0	0.1	3.079	0
–	–	–	–	–	0.008	–	–	–	–	–	–	–	–	3.088	1
–	–	–	–	–	–	–	–	–	–	–	256	–	–	3.083	2

Table 138: Hyperparameter ablation for Muon on 300m on 8x Chinchilla Data

$\eta_{\text{adam}}$	$\beta_1$	$\beta_2$	Decay(WSD)	$\epsilon$	$\eta$	Schedule	$g_{\text{norm}}$	$\eta_{\text{min}}$	$\beta_{\text{muon}}$	$\epsilon_{\text{muon}}$	BSZ	warmup	$\lambda$	Loss	Link
0.0024	0.8	0.98	0.8	1e-15	0.008	linear	1	0	0.98	1e-05	256	0	0.1	3.032	0
–	–	–	–	–	0.004	–	–	–	–	–	–	–	–	3.029	1
–	–	–	–	–	–	–	–	–	–	–	128	–	–	3.049	2

Table 139: Hyperparameter ablation for Muon on 520m on 2x Chinchilla Data

$\eta_{adam}$	$\beta_1$	$\beta_2$	Decay(WSD)	$\epsilon$	$\eta$	Schedule	$g_{norm}$	$\eta_{min}$	$\beta_{muon}$	$\epsilon_{muon}$	BSZ	warmup	$\lambda$	Loss	Link
0.0012	0.8	0.98	1	1e-25	0.004	linear	2	0	0.98	1e-05	128	0	0.1	3.002	0
–	–	–	–	–	0.008	–	–	–	–	–	–	–	–	3.008	1
–	–	–	–	–	–	–	–	–	–	–	256	–	–	3.009	2

Table 140: Hyperparameter ablation for Muon on 520m on 4x Chinchilla Data

$\eta_{adam}$	$\beta_1$	$\beta_2$	Decay(WSD)	$\epsilon$	$\eta$	Schedule	$g_{norm}$	$\eta_{min}$	$\beta_{muon}$	$\epsilon_{muon}$	BSZ	warmup	$\lambda$	Loss	Link
0.0024	0.8	0.98	1	1e-25	0.008	linear	2	0	0.98	1e-05	256	0	0.1	2.945	0
–	–	–	–	–	0.004	–	–	–	–	–	–	–	–	2.944	1
–	–	–	–	–	–	–	–	–	–	–	128	–	–	2.963	2

Table 141: Hyperparameter ablation for Muon on 520m on 8x Chinchilla Data

$\eta_{adam}$	$\beta_1$	$\beta_2$	Decay(WSD)	$\epsilon$	$\eta$	Schedule	$g_{norm}$	$\eta_{min}$	$\beta_{muon}$	$\epsilon_{muon}$	BSZ	warmup	$\lambda$	Loss	Link
0.0024	0.8	0.98	1	1e-25	0.008	linear	2	0	0.98	1e-05	256	0	0.1	2.906	0
–	–	–	–	–	0.004	–	–	–	–	–	–	–	–	2.900	1
–	–	–	–	–	–	–	–	–	–	–	128	–	–	2.930	2

## D.8 Sweeping Results for Scion

Table 142: Hyperparameter ablation for Scion on 300m on 2x Chinchilla Data

$\eta_{adam}$	$\beta_1$	Decay(WSD)	$\eta$	Schedule	$g_{norm}$	$\eta_{min}$	$\beta_{muon}$	$\epsilon_{scion}$	BSZ	warmup	$\lambda$	Loss	Link
0.0008	0.98	0.8	0.008	linear	2	0	0.95	1e-05	128	0	0.1	3.152	0
–	–	–	0.004	–	–	–	–	–	–	–	–	3.153	1
–	–	–	–	–	–	–	–	–	256	–	–	3.154	2

Table 143: Hyperparameter ablation for Scion on 300m on 4x Chinchilla Data

$\eta_{adam}$	$\beta_1$	Decay(WSD)	$\eta$	Schedule	$g_{norm}$	$\eta_{min}$	$\beta_{muon}$	$\epsilon_{scion}$	BSZ	warmup	$\lambda$	Loss	Link
0.0008	0.98	1	0.008	linear	2	0	0.95	1e-05	256	0	0.1	3.086	0
–	–	–	0.004	–	–	–	–	–	–	–	–	3.099	1
–	–	–	–	–	–	–	–	–	128	–	–	3.090	2

Table 144: Hyperparameter ablation for Scion on 300m on 8x Chinchilla Data

$\eta_{adam}$	$\beta_1$	Decay(WSD)	$\eta$	Schedule	$g_{norm}$	$\eta_{min}$	$\beta_{muon}$	$\epsilon_{scion}$	BSZ	warmup	$\lambda$	Loss	Link
0.0004	0.98	0.8	0.004	linear	2	0	0.95	1e-05	128	0	0.1	3.039	0
–	–	–	0.008	–	–	–	–	–	–	–	–	3.057	1
–	–	–	–	–	–	–	–	–	256	–	–	3.037	2

Table 145: Hyperparameter ablation for Scion on 520m on 2x Chinchilla Data

$\eta_{adam}$	$\beta_1$	Decay(WSD)	$\eta$	Schedule	$g_{norm}$	$\eta_{min}$	$\beta_{muon}$	$\epsilon_{scion}$	BSZ	warmup	$\lambda$	Loss	Link
0.0004	0.98	1	0.004	linear	2	0	0.95	1e-05	128	0	0.1	3.007	0
–	–	–	0.008	–	–	–	–	–	–	–	–	3.015	1
–	–	–	–	–	–	–	–	–	256	–	–	3.020	2

Table 146: Hyperparameter ablation for Scion on 520m on 4x Chinchilla Data

$\eta_{adam}$	$\beta_1$	Decay(WSD)	$\eta$	Schedule	$g_{norm}$	$\eta_{min}$	$\beta_{muon}$	$\epsilon_{scion}$	BSZ	warmup	$\lambda$	Loss	Link
0.0008	0.98	1	0.008	linear	2	0	0.95	1e-05	256	0	0.1	2.952	0
–	–	–	0.004	–	–	–	–	–	–	–	–	2.952	1
–	–	–	–	–	–	–	–	–	128	–	–	2.970	2

Table 147: Hyperparameter ablation for Scion on 520m on 8x Chinchilla Data

$\eta_{adam}$	$\beta_1$	Decay(WSD)	$\eta$	Schedule	$g_{norm}$	$\eta_{min}$	$\beta_{muon}$	$\epsilon_{scion}$	BSZ	warmup	$\lambda$	Loss	Link
0.0004	0.98	1	0.004	linear	2	0	0.95	1e-05	256	0	0.1	2.904	0
–	–	–	0.008	–	–	–	–	–	–	–	–	2.913	1
–	–	–	–	–	–	–	–	–	128	–	–	2.913	2

## D.9 Sweeping Results for Kron

Table 148: Hyperparameter ablation for Kron on 300m on 2x Chinchilla Data

$\beta_1$	blocksize	$\eta$	$g_{norm}$	$\eta_{min}$	NormGrad	Blocking	$Init_{pc}$	$\eta_{pc}$	$p_{pc}$	BSZ	$Step_{pc}$	warmup	$\lambda$	Loss	Link
0.95	256	0.001	1	0	True	True	1	0.2	0.1	128	2000	1000	0.7	3.151	0
–	–	0.0005	–	–	–	–	–	–	–	–	–	–	–	3.157	1

Table 149: Hyperparameter ablation for Kron on 300m on 4x Chinchilla Data

$\beta_1$	blocksize	$\eta$	$g_{norm}$	$\eta_{min}$	NormGrad	Blocking	$Init_{pc}$	$\eta_{pc}$	$p_{pc}$	BSZ	$Step_{pc}$	warmup	$\lambda$	Loss	Link
0.95	256	0.0005	1	0	True	True	1	0.2	0.1	128	2000	1000	0.7	3.083	0
–	–	0.001	–	–	–	–	–	–	–	–	–	–	–	3.090	1

Table 150: Hyperparameter ablation for Kron on 300m on 8x Chinchilla Data

$\beta_1$	blocksize	$\eta$	$g_{norm}$	$\eta_{min}$	NormGrad	Blocking	$Init_{pc}$	$\eta_{pc}$	$p_{pc}$	BSZ	$Step_{pc}$	warmup	$\lambda$	Loss	Link
0.95	256	0.0005	1	0	True	True	1	0.2	0.1	128	2000	1000	0.7	3.031	0
–	–	0.001	–	–	–	–	–	–	–	–	–	–	–	3.074	1

Table 151: Hyperparameter ablation for Kron on 520m on 2x Chinchilla Data

$\beta_1$	blocksize	$\eta$	$g_{\text{norm}}$	$\eta_{\text{min}}$	NormGrad	Blocking	$Init_{pc}$	$\eta_{pc}$	$p_{pc}$	BSZ	$Step_{pc}$	warmup	$\lambda$	Loss	Link
0.95	256	0.0005	1	0	True	True	1	0.2	0.1	128	2000	1000	0.5	3.009	0
–	–	0.001	–	–	–	–	–	–	–	–	–	–	–	3.009	1

Table 152: Hyperparameter ablation for Kron on 520m on 4x Chinchilla Data

$\beta_1$	blocksize	$\eta$	$g_{\text{norm}}$	$\eta_{\text{min}}$	NormGrad	Blocking	$Init_{pc}$	$\eta_{pc}$	$p_{pc}$	BSZ	$Step_{pc}$	warmup	$\lambda$	Loss	Link
0.95	256	0.0005	1	0	True	True	1	0.2	0.1	128	2000	1000	0.5	2.946	0
–	–	0.001	–	–	–	–	–	–	–	–	–	–	–	2.950	1

Table 153: Hyperparameter ablation for Kron on 520m on 8x Chinchilla Data

$\beta_1$	blocksize	$\eta$	$g_{\text{norm}}$	$\eta_{\text{min}}$	NormGrad	Blocking	$Init_{pc}$	$\eta_{pc}$	$p_{pc}$	BSZ	$Step_{pc}$	warmup	$\lambda$	Loss	Link
0.95	256	0.0005	1	0	True	True	1	0.2	0.1	128	2000	1000	0.5	2.900	0
–	–	0.001	–	–	–	–	–	–	–	–	–	–	–	2.909	1
–	–	–	–	–	–	–	–	–	–	256	–	–	–	2.902	2

## D.10 Sweeping Results for Soap

Table 154: Hyperparameter ablation for Soap on 300m on 2x Chinchilla Data

$\beta_1$	$\beta_2$	blocksize	$\epsilon$	$\eta$	$g_{\text{norm}}$	$\eta_{\text{min}}$	Blocking	$f_{pc}$	$\beta_{\text{shampoo}}$	BSZ	warmup	$\lambda$	Loss	Link
0.95	0.99	512	1e-10	0.008	1	0	True	10	0.9	128	1000	0.1	3.147	0
–	–	128	–	–	–	–	–	–	–	–	–	–	3.154	1
–	–	256	–	–	–	–	–	–	–	–	–	–	3.150	2
–	–	–	–	0.004	–	–	–	–	–	–	–	–	3.147	3
–	–	–	–	–	–	–	–	–	–	256	–	–	3.153	4

Table 155: Hyperparameter ablation for Soap on 300m on 4x Chinchilla Data

$\beta_1$	$\beta_2$	blocksize	$\epsilon$	$\eta$	$g_{\text{norm}}$	$\eta_{\text{min}}$	Blocking	$f_{pc}$	$\beta_{\text{shampoo}}$	BSZ	warmup	$\lambda$	Loss	Link
0.95	0.99	512	1e-10	0.008	1	0	True	10	0.9	256	1000	0.1	3.084	0
–	–	128	–	–	–	–	–	–	–	–	–	–	3.086	1
–	–	256	–	–	–	–	–	–	–	–	–	–	3.084	2
–	–	–	–	0.004	–	–	–	–	–	–	–	–	3.086	3
–	–	–	–	–	–	–	–	–	–	128	–	–	3.091	4

Table 156: Hyperparameter ablation for Soap on 300m on 8x Chinchilla Data

$\beta_1$	$\beta_2$	blocksize	$\epsilon$	$\eta$	$g_{\text{norm}}$	$\eta_{\text{min}}$	Blocking	$f_{pc}$	$\beta_{\text{shampoo}}$	BSZ	warmup	$\lambda$	Loss	Link
0.95	0.99	512	1e-10	0.008	1	0	True	10	0.9	256	1000	0.1	3.030	0
–	–	128	–	–	–	–	–	–	–	–	–	–	3.034	1
–	–	256	–	–	–	–	–	–	–	–	–	–	3.032	2
–	–	–	–	0.004	–	–	–	–	–	–	–	–	3.031	3
–	–	–	–	–	–	–	–	–	–	128	–	–	3.043	4

Table 157: Hyperparameter ablation for Soap on 520m on 2x Chinchilla Data

$\beta_1$	$\beta_2$	blocksize	$\epsilon$	$\eta$	$g_{\text{norm}}$	$\eta_{\text{min}}$	Blocking	$f_{pc}$	$\beta_{\text{shampoo}}$	BSZ	warmup	$\lambda$	Loss	Link
0.95	0.99	512	1e-10	0.008	1	0	True	10	0.95	256	1000	0.1	3.004	0
–	–	128	–	–	–	–	–	–	–	–	–	–	3.013	1
–	–	256	–	–	–	–	–	–	–	–	–	–	3.010	2
–	–	–	–	0.004	–	–	–	–	–	–	–	–	3.008	3
–	–	–	–	–	–	–	–	–	–	128	–	–	3.011	4

Table 158: Hyperparameter ablation for Soap on 520m on 4x Chinchilla Data

$\beta_1$	$\beta_2$	blocksize	$\epsilon$	$\eta$	$g_{\text{norm}}$	$\eta_{\text{min}}$	Blocking	$f_{pc}$	$\beta_{\text{shampoo}}$	BSZ	warmup	$\lambda$	Loss	Link
0.95	0.99	512	1e-10	0.004	1	0	True	10	0.95	256	1000	0.1	2.944	0
–	–	128	–	–	–	–	–	–	–	–	–	–	2.948	1
–	–	256	–	–	–	–	–	–	–	–	–	–	2.945	2
–	–	–	–	0.008	–	–	–	–	–	–	–	–	2.949	3
–	–	–	–	–	–	–	–	–	–	128	–	–	2.946	4

Table 159: Hyperparameter ablation for Soap on 520m on 8x Chinchilla Data

$\beta_1$	$\beta_2$	blocksize	$\epsilon$	$\eta$	$g_{\text{norm}}$	$\eta_{\text{min}}$	Blocking	$f_{pc}$	$\beta_{\text{shampoo}}$	BSZ	warmup	$\lambda$	Loss	Link
0.95	0.99	512	1e-10	0.004	1	0	True	10	0.95	256	1000	0.1	2.899	0
–	–	–	–	0.008	–	–	–	–	–	–	–	–	2.906	1

## E Hyperparameter Ablation in Phase III 1.2B experiments

We reported the results for the optimizers we swept in Phase III. The result is formulated in the same way as in Phase I.

## E.1 Sweeping Results for AdamW

Table 160: Hyperparameter ablation for AdamW on 1.2b on 1x Chinchilla Data

$\beta_1$	$\beta_2$	$\epsilon$	$\eta$	$g_{\text{norm}}$	$\eta_{\text{min}}$	BSZ	warmup	$\lambda$	Loss	Link
0.9	0.98	1e-10	0.002	2	0	256	2000	0.2	2.905	0
0.95	–	–	–	–	–	–	–	–	2.905	1
0.98	–	–	–	–	–	–	–	–	2.909	2
–	0.9	–	–	–	–	–	–	–	2.914	3
–	0.95	–	–	–	–	–	–	–	2.909	4
–	–	1e-25	–	–	–	–	–	–	2.907	5
–	–	1e-20	–	–	–	–	–	–	2.907	6
–	–	1e-15	–	–	–	–	–	–	2.907	7
–	–	–	0.004	–	–	–	–	–	2.916	8
–	–	–	0.008	–	–	–	–	–	7.347	9
–	–	–	–	0	–	–	–	–	2.909	10
–	–	–	–	1.0	–	–	–	–	2.908	11
–	–	–	–	–	–	128	–	–	2.904	12
–	–	–	–	–	–	512	–	–	2.928	13
–	–	–	–	–	–	1024	–	–	2.985	14
–	–	–	–	–	–	–	500	–	2.917	15
–	–	–	–	–	–	–	1000	–	2.910	16
–	–	–	–	–	–	–	4000	–	2.912	17
–	–	–	–	–	–	–	–	0	2.946	18
–	–	–	–	–	–	–	–	0.1	2.916	19

Table 161: Hyperparameter ablation for AdamW on 1.2b on 2x Chinchilla Data

$\beta_1$	$\beta_2$	$\epsilon$	$\eta$	$g_{\text{norm}}$	$\eta_{\text{min}}$	BSZ	warmup	$\lambda$	Loss	Link
0.9	0.98	1e-10	0.002	1	0	256	1000	0.2	2.836	0

Table 162: Hyperparameter ablation for AdamW on 1.2b on 4x Chinchilla Data

$\beta_1$	$\beta_2$	$\epsilon$	$\eta$	$g_{\text{norm}}$	$\eta_{\text{min}}$	BSZ	warmup	$\lambda$	Loss	Link
0.9	0.98	1e-10	0.002	1	0	256	1000	0.2	2.787	0

Table 163: Hyperparameter ablation for AdamW on 1.2b on 8x Chinchilla Data

$\beta_1$	$\beta_2$	$\epsilon$	$\eta$	$g_{\text{norm}}$	$\eta_{\text{min}}$	BSZ	warmup	$\lambda$	Loss	Link
0.9	0.98	1e-10	0.002	1	0	256	1000	0.2	2.752	0



## E.2 Sweeping Results for NAdamW

Table 164: Hyperparameter ablation for NAdamW on 1.2b on 1x Chinchilla Data

$\beta_1$	$\beta_2$	$\epsilon$	$\eta$	$g_{\text{norm}}$	$\eta_{\text{min}}$	BSZ	warmup	$\lambda$	Loss	Link
0.98	0.98	1e-10	0.004	1	0	256	4000	0.1	2.902	0

Table 165: Hyperparameter ablation for NAdamW on 1.2b on 2x Chinchilla Data

$\beta_1$	$\beta_2$	$\epsilon$	$\eta$	$g_{\text{norm}}$	$\eta_{\text{min}}$	BSZ	warmup	$\lambda$	Loss	Link
0.98	0.98	1e-10	0.004	1	0	256	4000	0.1	2.833	0

Table 166: Hyperparameter ablation for NAdamW on 1.2b on 4x Chinchilla Data

$\beta_1$	$\beta_2$	$\epsilon$	$\eta$	$g_{\text{norm}}$	$\eta_{\text{min}}$	BSZ	warmup	$\lambda$	Loss	Link
0.98	0.98	1e-10	0.004	1	0	256	4000	0.1	2.785	0

Table 167: Hyperparameter ablation for NAdamW on 1.2b on 8x Chinchilla Data

$\beta_1$	$\beta_2$	$\epsilon$	$\eta$	$g_{\text{norm}}$	$\eta_{\text{min}}$	BSZ	warmup	$\lambda$	Loss	Link
0.98	0.98	1e-10	0.004	1	0	256	4000	0.1	2.749	0

## E.3 Sweeping Results for Soap

Table 168: Hyperparameter ablation for Soap on 1.2b on 1x Chinchilla Data

$\beta_1$	$\beta_2$	blocksize	$\epsilon$	$\eta$	$g_{\text{norm}}$	$\eta_{\text{min}}$	Blocking	$f_{\text{pc}}$	$\beta_{\text{shampoo}}$	BSZ	warmup	$\lambda$	Loss	Link
0.95	0.99	512	1e-10	0.004	1	0	True	10	0.9	256	1000	0.1	2.940	0

Table 169: Hyperparameter ablation for Soap on 1.2b on 2x Chinchilla Data

$\beta_1$	$\beta_2$	blocksize	$\epsilon$	$\eta$	$g_{\text{norm}}$	$\eta_{\text{min}}$	Blocking	$f_{\text{pc}}$	$\beta_{\text{shampoo}}$	BSZ	warmup	$\lambda$	Loss	Link
0.95	0.99	512	1e-10	0.004	1	0	True	10	0.9	256	1000	0.1	2.829	0

Table 170: Hyperparameter ablation for Soap on 1.2b on 4x Chinchilla Data

$\beta_1$	$\beta_2$	blocksize	$\epsilon$	$\eta$	$g_{\text{norm}}$	$\eta_{\text{min}}$	Blocking	$f_{\text{pc}}$	$\beta_{\text{shampoo}}$	BSZ	warmup	$\lambda$	Loss	Link
0.95	0.99	512	1e-10	0.004	1	0	True	10	0.9	256	1000	0.1	2.783	0

Table 171: Hyperparameter ablation for Soap on 1.2b on 8x Chinchilla Data

$\beta_1$	$\beta_2$	blocksize	$\epsilon$	$\eta$	$g_{\text{norm}}$	$\eta_{\text{min}}$	Blocking	$f_{\text{pc}}$	$\beta_{\text{shampoo}}$	BSZ	warmup	$\lambda$	Loss	Link
0.95	0.99	512	1e-10	0.004	1	0	True	10	0.9	256	1000	0.1	2.749	0

## E.4 Sweeping Results for Muon

Table 172: Hyperparameter ablation for Muon on 1.2b on 1x Chinchilla Data

$\eta_{\text{adam}}$	$\beta_1$	$\beta_2$	Decay(WSD)	$\epsilon$	$\eta$	Schedule	$g_{\text{norm}}$	$\eta_{\text{min}}$	$\beta_{\text{muon}}$	$\epsilon_{\text{muon}}$	BSZ	warmup	$\lambda$	Loss	Link
0.0012	0.8	0.98	1	1e-15	0.004	linear	2	0	0.98	1e-05	256	0	0.1	2.891	0
–	–	–	–	–	0.008	–	–	–	–	–	–	–	–	2.886	1

Table 173: Hyperparameter ablation for Muon on 1.2b on 2x Chinchilla Data

$\eta_{\text{adam}}$	$\beta_1$	$\beta_2$	Decay(WSD)	$\epsilon$	$\eta$	Schedule	$g_{\text{norm}}$	$\eta_{\text{min}}$	$\beta_{\text{muon}}$	$\epsilon_{\text{muon}}$	BSZ	warmup	$\lambda$	Loss	Link
0.0012	0.8	0.98	1	1e-15	0.004	linear	2	0	0.98	1e-05	256	0	0.1	2.827	0
–	–	–	–	–	0.008	–	–	–	–	–	–	–	–	2.833	1

Table 174: Hyperparameter ablation for Muon on 1.2b on 4x Chinchilla Data

$\eta_{\text{adam}}$	$\beta_1$	$\beta_2$	Decay(WSD)	$\epsilon$	$\eta$	Schedule	$g_{\text{norm}}$	$\eta_{\text{min}}$	$\beta_{\text{muon}}$	$\epsilon_{\text{muon}}$	BSZ	warmup	$\lambda$	Loss	Link
0.0012	0.8	0.98	1	1e-15	0.004	linear	2	0	0.98	1e-05	256	0	0.1	2.780	0
–	–	–	–	–	0.008	–	–	–	–	–	–	–	–	2.793	1

Table 175: Hyperparameter ablation for Muon on 1.2b on 8x Chinchilla Data

$\eta_{\text{adam}}$	$\beta_1$	$\beta_2$	Decay(WSD)	$\epsilon$	$\eta$	Schedule	$g_{\text{norm}}$	$\eta_{\text{min}}$	$\beta_{\text{muon}}$	$\epsilon_{\text{muon}}$	BSZ	warmup	$\lambda$	Loss	Link
0.0012	0.8	0.98	1	1e-15	0.004	linear	2	0	0.98	1e-05	256	0	0.1	2.748	0

## F Hyperparameter Ablation in Phase III 16x Chinchilla experiments

### F.1 Sweeping Results for AdamW

Table 176: Hyperparameter ablation for AdamW on 130m on 16x Chinchilla Data

$\beta_1$	$\beta_2$	$\epsilon$	$\eta$	$g_{\text{norm}}$	$\eta_{\text{min}}$	BSZ	warmup	$\lambda$	Loss	Link
0.9	0.98	1e-10	0.008	1	0	256	1000	0.1	3.207	0

Table 177: Hyperparameter ablation for AdamW on 300m on 16x Chinchilla Data

$\beta_1$	$\beta_2$	$\epsilon$	$\eta$	$g_{\text{norm}}$	$\eta_{\text{min}}$	BSZ	warmup	$\lambda$	Loss	Link
0.9	0.98	1e-10	0.004	2	0	256	2000	0.1	3.001	0

## F.2 Sweeping Results for NAdamW

Table 178: Hyperparameter ablation for NAdamW on 130m on 16x Chinchilla Data

$\beta_1$	$\beta_2$	$\epsilon$	$\eta$	$g_{\text{norm}}$	$\eta_{\text{min}}$	BSZ	warmup	$\lambda$	Loss	Link
0.95	0.98	1e-10	0.008	1	0	256	2000	0.1	3.200	0

Table 179: Hyperparameter ablation for NAdamW on 300m on 16x Chinchilla Data

$\beta_1$	$\beta_2$	$\epsilon$	$\eta$	$g_{\text{norm}}$	$\eta_{\text{min}}$	BSZ	warmup	$\lambda$	Loss	Link
0.95	0.98	1e-10	0.004	1	0	256	2000	0.1	2.998	0

## F.3 Sweeping Results for Soap

Table 180: Hyperparameter ablation for Soap on 130m on 16x Chinchilla Data

$\beta_1$	$\beta_2$	blocksize	$\epsilon$	$\eta$	$g_{\text{norm}}$	$\eta_{\text{min}}$	Blocking	$f_{\text{pc}}$	$\beta_{\text{shampoo}}$	BSZ	warmup	$\lambda$	Loss	Link
0.95	0.99	512	1e-10	0.008	1	0	True	10	0.98	256	1000	0.1	3.191	0

Table 181: Hyperparameter ablation for Soap on 300m on 16x Chinchilla Data

$\beta_1$	$\beta_2$	blocksize	$\epsilon$	$\eta$	$g_{\text{norm}}$	$\eta_{\text{min}}$	Blocking	$f_{\text{pc}}$	$\beta_{\text{shampoo}}$	BSZ	warmup	$\lambda$	Loss	Link
0.95	0.99	512	1e-10	0.004	1	0	True	10	0.9	256	1000	0.1	2.990	0

## F.4 Sweeping Results for Muon

Table 182: Hyperparameter ablation for Muon on 130m on 16x Chinchilla Data

$\eta_{\text{adam}}$	$\beta_1$	$\beta_2$	Decay(WSD)	$\epsilon$	$\eta$	Schedule	$g_{\text{norm}}$	$\eta_{\text{min}}$	$\beta_{\text{muon}}$	$\epsilon_{\text{muon}}$	BSZ	warmup	$\lambda$	Loss	Link
0.0012	0.8	0.98	1	1e-25	0.004	linear	1	0	0.98	1e-05	128	0	0.1	3.192	0
–	–	–	–	–	0.008	–	–	–	–	–	–	–	–	3.202	1

Table 183: Hyperparameter ablation for Muon on 300m on 16x Chinchilla Data

$\eta_{\text{adam}}$	$\beta_1$	$\beta_2$	Decay(WSD)	$\epsilon$	$\eta$	Schedule	$g_{\text{norm}}$	$\eta_{\text{min}}$	$\beta_{\text{muon}}$	$\epsilon_{\text{muon}}$	BSZ	warmup	$\lambda$	Loss	Link
0.0012	0.8	0.98	0.8	1e-15	0.004	linear	1	0	0.98	1e-05	256	0	0.1	2.991	0

## G Comparison with Prior Work

This paper benchmarks the performance of 11 proposed optimizers and show vastly different speed-up ratio than prior works reported. In this section, we will compare the setup of our experiments with prior works with the hope of understanding the difference.

1. Sophia Liu et al. [2024a] ( $2\times$ ) This work utilizes a small peak learning rate of learning rate smaller than  $6e-4$  (similar to the one shown in Figure 1 Top Left). The reason for the small peak learning rate is likely 2-fold: (i) the authors are training on a pretraining dataset PILE with lower quality compared to the current pretraining dataset and (ii) in the implementation of Levanter that the authors used, the data shuffling is not completely random and instead is correlated on every compute node. Upon reproducing the results, we note that this difference can significantly impact the stability of the training process and a complete random shuffling is crucial for the usage of a large learning rate.
2. MARS Yuan et al. [2025] ( $2\times$ ) This papers considers a similar setup as Sophia and uses a similar AdamW baseline. We note that in the first version of the paper, the authors also reported that increasing the learning rate of AdamW to  $3e-3$  can significantly improve the performance of AdamW (see Figure 6 of Yuan et al. [2025] arxiv version 1).
3. Soap Vyas et al. [2025] ( $1.4\times$ ) The actual speedup of Soap on 300M and 520M models are 1.2 to  $1.3\times$ , which is only slightly lower than the claimed  $1.4\times$  speedup. We note that our implementation of Soap is slightly different from the one used in the paper that we performs blocking of weight in order to reduce the memory footprint and further uses bfloat16 for the momentum in the 1.2B experiments. Both modifications may lead to slightly lower step-wise performance.
4. Muon Jordan et al. [2024], Liu et al. [2025a] ( $2\times$ ) The speedup of Muon reported in different works are vastly different. In the original Nanogpt speedrun, Muon achieves  $1.3\times$  speedup over AdamW. Later the reproduction of Kimi reported a much higher speedup of  $2\times$ . We note that the Kimi version utilizes a notably low learning rate for AdamW ( $8e-4$  to  $9e-4$  for model between 400M to 1.5B) in the scaling experiments. We also note that the smaller learning rate is important for hyperparameter transfer from AdamW when roughly matching the update norm and AdamW can perform better higher learning rate in our experiments. Further, their comparison of Muon and AdamW on the MoE experiments compare two models with not fully decayed learning rate and this may significantly favors Muon, as shown in Figure 5. It is later shown independently in the work of Essential AI AI et al. [2025] that Muon’s token efficiency compared to AdamW is only 1.1 to  $1.2\times$ .
5. Cautious Liang et al. [2025], Block-wise Learning Rate Adam Wang et al. [2025], FOCUS Liu et al. [2025c] report  $2\times$  speedup over AdamW. These papers use a similar baseline as Sophia and MARS.
6. SWAN Ma et al. [2025] and DION Ahn et al. [2025] report  $2-3\times$  speedup over AdamW. The comparison between these two optimizers and AdamW is carried out on a smaller than  $1\times$  Chinchilla regime, where the speed-up of matrix-based optimizer may be larger. We also note that we didn’t consider the communication cost, which is the main focus of DION.
7. SPlus Frans et al. [2025] reports a  $2\times$  speedup over AdamW. This work considers an atypical setup where the model is trained with a constant learning rate.