SCALING LAW WITH LEARNING RATE ANNEALING

Anonymous authors

000

001 002 003

004

005 006 007

008 009

010

011

013

014

015

016

017

018

019

021

023

025

026

027

028 029

031

Paper under double-blind review

ABSTRACT

We find that the cross-entropy loss curves of neural language models empirically adhere to a scaling law with learning rate (LR) annealing over training steps:

$$L(s) = L_0 + A \cdot S_1^{-\alpha} - C \cdot S_2$$

where L(s) is the validation loss at step s, S_1 is the area under the LR curve, S_2 is the LR annealing area, and L_0 , A, C, α are constant parameters. This formulation accounts for two main effects: (1) power-law scaling over data size, and (2) the additional loss reduction during LR annealing. Unlike previous studies that only fit losses at final steps, our formulation captures the entire training curve, allowing for parameter fitting using losses from any training step. Applying the scaling law with LR annealing and fitting only one or two training curves, we can accurately predict the loss at any given step under any learning rate scheduler (LRS). This approach significantly reduces computational cost in formulating scaling laws while providing more accuracy and expressiveness. Extensive experiments demonstrate that our findings hold across a range of hyper-parameters and model architectures and can extend to scaling effect of model sizes. Moreover, our formulation provides accurate theoretical insights into empirical results observed in numerous previous studies, particularly those focusing on LR schedule and annealing. We believe that this work is promising to enhance the understanding of LLM training dynamics while democratizing scaling laws, and it is helpful to guide both research and industrial participants in refining training strategies for further LLMs.

030 1 INTRODUCTION

In recent years, large language models (LLMs) have garnered significant academic and industrial attention (Brown et al., 2020; Touvron et al., 2023). The scaling law suggests that the validation loss of language models follow a power-law pattern as model and data sizes increase (Hestness et al., 2017; Kaplan et al., 2020; Henighan et al., 2020). This law provides a powerful framework for forecasting LLM performances before large scale training by fitting losses at smaller scales (OpenAI, 2023; DeepSeek-AI, 2024; Dubey et al., 2024). Numerous studies have explored on the formulation to model the scaling effect of LLMs under various different settings (Bahri et al., 2021; Hernandez et al., 2021; Caballero et al., 2022; Michaud et al., 2023; Muennighoff et al., 2023).

However, typical scaling law formulations focus only on the final loss at the end of training (Hoff-040 mann et al., 2022). Specifically, previous approaches generally rely on a set of training runs and 041 fit the scaling law curve solely on the final loss from each run, while ignoring middle losses during 042 training which do not follow traditional scaling laws. This approach underutilizes the training com-043 pute and fails to capture the training dynamics within each run. Further, the application of scaling 044 laws in LLM developments is limited since the loss curve through the whole training process is 045 not modeled. An expressive formulation that models full loss curves enables prediction of future 046 training dynamics and also offers insights on understanding the learning process of LLMs. 047

In this study, we propose a scaling law that models the full loss curve within a complete LLM training run. Specifically, we dive deeper into the training dynamics during LR annealing, and incorporate a LR annealing factor into the traditional scaling law formula to formulate the process. This design is motivated by the observed correlation between LRS and loss curves, where loss gradually decreases as we consume more training steps ¹ and then sharply declines when the LR

¹In this paper, we use training steps to quantify the amount of consumed data, as they are typically proportional, with data amount calculated as training steps multiplied by batch size.

055

056

058

059

060

061

062 063

064 065

067

068

069

071

072

073 074

075

076

077 078 079

080 081



Figure 1: Visualization of S_1 and S_2 at the 20-th step of a cosine LR scheduler. S_1 is the forward area, i.e., sum of red grid areas; S_2 is the decayed annealing area, i.e., weighted sum of blue grid areas, where lighter shades indicate smaller weights. Both S_1 and S_2 contribute to loss reduction, and balancing their values is crucial for achieving the lowest possible final loss.

undergoes significant annealing (Loshchilov & Hutter, 2016; Smith et al., 2018; Ibrahim et al., 2024; Hu et al., 2024). We propose that the model's validation loss L(s) at step s is determined by two main factors: the forward area S_1 under the LR curve and the degree of LR annealing S_2 :

$$L(s) = L_0 + A \cdot S_1^{-\alpha} - C \cdot S_2,$$

$$S_1 = \sum_{i=1}^s \eta_i, \qquad S_2 = \sum_{i=1}^s \sum_{k=1}^i (\eta_{k-1} - \eta_k) \cdot \lambda^{i-k},$$
(1)

where η_i is the learning rate at step *i*, and λ is a hyper-parameter representing the decay factor for LR annealing momentum (see Sec. 3 in detail), which typically ranges from 0.99 to 0.999. L_0 , A, C, α are undetermined positive constants. S_1 is also known as the summed learning rate (Kaplan et al., 2020), and S_2 represents the LR annealing area. A visualization of S_1 and S_2 is provided as Fig. 1.

Eq. 1 describes how loss changes for each step in a full loss curve during training. In Eq. 1, the term $L_0 + A \cdot S_1^{-\alpha}$ represents a rectified scaling law that captures the expected loss decreases as a power-law function of the number of training steps. The new term $-C \cdot S_2$ accounts for the further loss drop due to learning rate annealing. Remarkably, this simple formulation accurately describes the validation loss at any training step across various LRS and even allows us to predict the loss curve for unseen LRS. For example, we can fit Eq. 1 to the full loss curve of constant and cosine LRS with 20K total steps (Fig. 2), and then predict the full loss curve for various unseen LRS with longer total steps (e.g. 60K) (Fig. 3).

We validate our proposed equation through extensive experiments and find that: (1) Our formulation performs consistently well across various hyper-parameters and model architectures; (2) Eq. 1 can be extended to incorporate other scaling factors, such as model sizes; (3) Our proposed equation accurately fits the loss curves of open-sourced models; (4) Our formulation can be used to verify and explain numerous previous findings regarding LR annealing and scheduling.

100 In Sec. 3, we derive the scaling law formulation with LR annealing and elucidate the potential theory 101 underpinning our formulation. Extensive experiments are conducted to validate the formulation. In 102 Sec. 4, we apply our formulation to verify and explain the empirical results from various previous 103 studies. Our approach offers theoretical insights into the crux of loss drop, LR schedule, and LR 104 annealing, enabling LLM participants to better understand training dynamics of LLM and select op-105 timal training recipes in advance. In Sec. 5, we compare our approach to typical scaling law formula, such as the Chinchilla scaling law (Hoffmann et al., 2022). We show that our formulation is more general and requires significantly less compute to fit, which greatly democratizes the development 107 of LLMs and scaling laws.





162 2 PRELIMINARY 163

164 Scaling Laws. Cross-entropy loss of language models on the validation set is a reliable indicator of LLMs' performance on downstream tasks (Caballero et al., 2022; Du et al., 2024). Kaplan et al. 165 (2020) empirically discovered a power-law relationship between validation loss L and three factors: 166 model size N, dataset size D, and training compute. As an application of scaling law, Hoffmann 167 et al. (2022) developed Chinchilla, a compute-optimal LLM, by balancing model size and dataset 168 size. They used a simplified and intuitive scaling law equation: $L(D, N) = L_0 + A \cdot D^{-\alpha} + B \cdot N^{-\beta}$, where L_0 , A, B, α , β are positive constants. Traditional scaling law formulations fit only the loss 170 at the final training step, while ignoring losses from other steps. Collecting a new loss value of data 171 size requires launching a another training run with the same LRS, which is resource-intensive. 172

Learning Rate Annealing. Learning rate annealing is a widely-used technique in training neural 173 networks, where the learning rate is progressively reduced from a maximum to a minimum value 174 following a pre-defined LRS. Various LRS schemes have been proposed to improve the performance 175 and stability of model training (Loshchilov & Hutter, 2016). For example, the popular cosine LRS 176 reduces the LR in a cosine-like pattern over full training steps. WSD LRS (Hu et al., 2024) keeps a 177 constant LR for the majority of training, and applies annealing only in the final (e.g. $10\% \sim 20\%$) 178 steps. In LLM training, it has been widely observed that a more pronounced decrease in the learning 179 rate often results in a more precipitous drop in the validation loss.

180 181

182

183

207 208

OBSERVATIONS AND EXPERIMENTS 3

In this section, we elaborate the origin, the intuition, and the experimental basis behind Eq. 1. We then validate our formula through extensive experiments. 184

185 3.1 SIMILARITY BETWEEN LEARNING RATE, GRADIENT NORM, AND LOSS

187 The first key observation is that the shapes of 188 LR curve, gradient norm curve, and validation 189 loss curve are quite similar across various LRS 190 when training LLMs (Fig. 4). This suggests an 191 implicit connection between learning rate and loss, where gradient norm could be the bridge. 192

193 Scaling Laws for Constant LRS. A constant 194 LRS is a special LRS, where every training 195 step can be viewed as an endpoint of the LRS. 196 Notably, the Chinchilla scaling law (Hoffmann et al., 2022) exactly fits losses of last steps, i.e., 197 LRS endpoints. Therefore, it is expected that the validation loss of all steps under a constant 199 LRS adheres to the Chinchilla scaling law, i.e., 200 a power-law over training step s. 201



Figure 4: The shapes of LR (top), gradient norm (medium), and validation loss (bottom) curves exhibit high similarity across various LRS (labeled as different colors).

202 Extra Loss Changes in LR Annealing. Unlike a constant LRS, LR annealing (or re-warmup) 203

brings significant local changes in the loss (see Fig. 4), causing the full loss curve to deviate from 204 the traditional power-law formulation that consider only the training steps s. We hypothesis that 205 such loss changes can be captured by an additional LR (η) related term, i.e., 206

$$L(s) = L_0 + A \cdot s^{-\alpha} - f(\eta),$$
(2)

where the first two terms follow traditional scaling laws, while the last term denotes the extra loss 209 change brought by LR annealing. Recall the similarity between learning rate and loss curves, we 210 can form an initial guess for $f(\eta)$ as $f(\eta) = C \cdot \eta$, where C is a positive constant. 211

212 **Training Discount in Annealing.** The form of Eq. 2 is still imperfect. Note that the gradient 213 norm $\|\mathbf{g}\|$ decreases almost proportionally with LR during the annealing process (shown in Fig. 4). Thus, the amount of parameter movement (approximately $\eta \cdot \|\mathbf{g}\|$ per step) in the LR annealing 214 stage declines at an almost quadratic rate compared to stages before annealing. As the parameter 215 movement become smaller, the change in loss also slows down accordingly. Therefore, the loss drop



(a) Different delay steps in the annealing process associated with different annealing steps (0.1K, 0.5K, 1K and 2K).

229

230 231

232 233 234

235

240

241 242

264 265 266



(b) Different delay steps in the re-warmup process associated with different re-warmup steps (0.1K, 0.5K, 1K and 2K).

Figure 5: The delay phenomenon between the LR and validation loss curves. This phenomenon suggests that LR annealing (re-warmup) has momentum.

brought by the power law term (i.e., the first two terms in Eq. 2) should also diminish during LR annealing. This consideration leads to an improved equation:

$$L(s) = L_0 + A \cdot S_1^{-\alpha} - f(\eta), \qquad S_1 = \sum_{i=1}^{s} \eta_i,$$
(3)

where S_1 is the forward area, i.e., the area under the LR curve (as visualized in Fig.1), which could be approximately interpreted as the total amount of parameter updates.

3.2 LR ANNEALING MOMENTUM 243

244 Another key observation is that LR annealing has momentum. To refine the formulation of $f(\eta)$, 245 we design a special LRS where the LR decreases linearly from η_{max} to η_{min} and then increases. 246 The increasing stage always has a fixed slope, reaching the maximum value in 5K steps, while the 247 slope of the decreasing stage is varied, with durations of 0.1K, 0.5K, 1K, and 2K. Symmetrically, we 248 design another LRS where the LR increases linearly from η_{min} to η_{max} and then decreases. Fig. 5 249 shows the corresponding LR and loss curves.

250 We observe a **delay phenomenon** between the LR and the validation loss. Firstly, the turning point 251 of the validation loss curve consistently lags behind the turning point of the LR curve, indicating that the validation loss continuous along its previous trajectory for some steps even after the LR 253 changes direction. Secondly, the steeper the slope of the LR annealing (or re-warmup), the more 254 pronounced the delay phenomenon becomes. Thirdly, given the same LR slope, the left figure 255 (where LR decreases then increases) consistently shows a longer delay compared to the right figure 256 (where LR increases then decreases).

257 Interestingly, this phenomenon closely resembles the physical experiment of a small ball rolling 258 down a slope. The steeper the slope, the faster the ball accelerates. When the ball lands, the accumu-259 lated momentum causes the ball to slide further. Inspired by this delay phenomenon, we hypothesize 260 that $f(\eta)$, the loss reduction induced by LR annealing, has cumulative historical formation so that 261 the past change of learning rate will affect the following loss curve for a few steps. In summary, 262 *learning rate annealing exhibits momentum.* To capture this, we define $f(\eta) = C \cdot S_2$, where S_2 is calculated as: 263

$$m_{i} = \lambda \cdot m_{i-1} + (\eta_{i-1} - \eta_{i}),$$

$$S_{2} = \sum_{i=1}^{s} m_{i} = \sum_{i=1}^{s} \sum_{k=1}^{i} (\eta_{k-1} - \eta_{k}) \cdot \lambda^{i-k},$$
(4)

where
$$m_i$$
 is the LR annealing momentum at step i ($m_1 = 0$), and $\Delta \eta = \eta_{i-1} - \eta_i$ denotes the Li
annealing amount at step i . λ is the decay factor that signifies how much historical information is

R is retained. We find that λ values between 0.99 and 0.999 generally works well. In contrast, $\lambda = 0$ implies no momentum effect, reducing $f(\eta)$ to $C \cdot \eta_s$, which degenerate to the initial form mentioned above. Note that S_2 applies not only to LR annealing ($S_2 > 0$), but also to LR re-warmup ($S_2 < 0$). This means that our equation is applicable to scenarios like continual pre-training, where LR rewarmup plays an important role in improving outcomes. Fig. 1 presents a visualization of S_2 .

275 3.3 FINAL FORMULATION

274

We formally present our formulation for the scaling law with LR annealing:

Scaling Law with LR Annealing. Given the same training and validation dataset, the same model size, the same training hyper-parameters such as warmup steps, **max learning rate** η_{max} and batch size, the language modeling loss at training step s empirically follows the equation $L(s) = L_0 + A \cdot S_1^{-\alpha} - C \cdot S_2$, where S_1 and S_2 are defined in Eq. 1. L_0 , A, C, α are positive constants.

Our formulation describes the loss of each training step across different LRS. It allows fitting based on a simpler LRS with shorter training steps and enables the prediction of validation losses for more complex LRS with longer training steps. Notably, loss curves with different max learning rates have different values of L_0 , A, C, α , and our scaling law does not fit divergent and collapsed loss curves (e.g., overly large LR). We also discuss some possible corner cases (i.e., $\eta = 0$) in Appendix H.3.

Loss Surface as a Slide. To better understand our formulation, we view the loss surface of language models as a slide in Fig. 6. The optimization process can be seen as sliding down the slide according to the power-law scaling (orange line), while oscillating on the inner wall (blue dashed line). When the LR anneals (red line), the amplitude of the oscillation decreases, resulting in a reduction in loss.

291 **Balance between** S_1 and S_2 . Note that in Eq.1, 292 $\frac{\partial L}{\partial S_1} < 0$ and $\frac{\partial L}{\partial S_2} < 0$ always hold, indicating that 293 increases in both S_1 and S_2 help to reduce the loss. However, as shown intuitively in Fig. 1, there ex-295 ists delicate balance between S_1 and S_2 . When LR 296 begins to anneal and S_2 starts to increase, the for-297 ward area S_1 of subsequent steps starts to diminish 298 instead. Our equation aptly describes this delicate 299 balance. In Sec. 4, we elaborate this topic in detail.

301 3.4 EXPERIMENTS

300

302 LR Warmup. LR warmup is important for train-303 ing LLMs. During the warmup stage, neural net-304 works are prone to random optimization, resulting 305 in unpredictable outcomes (Hestness et al., 2017). 306 Various studies, along with our own pilot experi-307 ments (Appendix A), show that LR warmup signif-308 icantly accelerates model convergence. High gra-309 dient norms are usually observed during the LR warmup stage, especially in the initial steps of train-310 ing (see Fig. 4). This indicates that model param-311 eters undergo substantial updates during this stage. 312



Figure 6: Loss surface of language models as a *slide* after simplification. Optimization direction could be decomposed into two directions: power-law scaling direction (S_1 , sliding down) and annealing direction (S_2 , inner height of the slide).

Therefore, in all our experiments, we linearly warmup LR to reach η_{max} and compute S_1 and S_2 assuming a constant LR value η_{max} in the warmup stage.

Experimental Setups and Fitting Details. We use standard experimental setups for LLM pretraining. To verify the robustness of our formulation across different experimental settings, we have five distinct experimental setups (see Appendix C). We adopt $\lambda = 0.999$ in our all experiments. We follow the fitting approach of Hoffmann et al. (2022) to obtain parameters in our equation (see Appendix B for more details).

Fitting and Prediction Results. We fit Eq.1 on the loss curves under constant and cosine LRS with 20K total steps (see Fig. 2), and then predict the full loss curves under several unseen LRS with 60K total steps (see Fig. 3). The results show an almost perfect fit, achieving a coefficient of determination (R^2) greater than 0.999. This underscores the robust capability of our equation to accurately fit loss curves across diverse LRS using a single parameter tuple.



(a) The loss drop brought by LR annealing for different model size N. Dashed lines represent trends over steps. The loss drop brought by LR annealing scales with data size and model size.

(b) Curve fitting and prediction on cosine LRS (60K steps to $\eta_{min} = 0.1 \cdot \eta_{max}$) of different model sizes using our *N*-extended scaling law. Results for N=1214.4M are predicted.

Figure 7: The loss drop brought by LR annealing (left) and the *N*-extended full loss curve fitting and prediction (right).

The prediction results in Fig. 3 indicate that our formulation is broadly applicable and generalizes robustly across four unseen LRS, with a mean prediction error as low as 0.2%. Moreover, our equation can accurately predict losses even for complex LRS that include multiple LR re-warmup stages (Fig. 3d), despite that the loss curves used for fitting do not contain any LR re-warmup stages.

347 Extensive Experiments on Different Setups. To demonstrate the broad applicability of our pro-348 posed equation, we conduct additional fitting and prediction experiments using various setups. (1) 349 We use an alternative set of training hyper-parameters (Appendix D.1); (2) We test our equation on the Mixture of Experts (MoE) architecture (Appendix D.2); (3) We apply our equation to predict 350 loss curves for a much longer training run involving a 1.7B parameter model trained on 1.4T tokens 351 (Appendix D.3). (4) We fit the loss curves of open-sourced models, including BLOOM-176B trained 352 on 300B tokens (BigScience, 2022) and OLMo-1B trained on 2T tokens (Groeneveld et al., 2024) 353 (Appendix D.4). All experiments produce excellent results, indicating that our equation is effec-354 tive across diverse experimental setups, including different training hyper-parameters, architecture, 355 model sizes, and dataset scales. We also present the ablation studies on S_1 and S_2 in Appendix D.5, 356 which shows each component in our formulation is important and indispensable. 357

358 359

367

370

3.5 EXTENSION TO MODEL SIZE SCALING

Loss Drop During Annealing Scales with Model Size N. We explore the effect of model size Non the loss drop during the annealing stage. Specifically, we compare the final losses obtained with a constant LRS and a WSD LRS (10% cosine annealing to $\eta_{min} = 0$) to estimate the loss drop due to LR annealing. We conduct this experiment on different total steps and different model sizes. The experimental results are shown in Fig. 7a. It suggest that the loss drop from LR annealing scales with both annealing steps and model sizes. This implies that the annealing area S_2 in our equation should also increase as the model size N increases. We suppose there is a simple relationship of $S_2 \propto N^{\gamma}$ where γ is a positive constant.

Model Size Scaling. Building on the experiments and analysis above, we extend our proposed Eq.1 to incorporate model size scaling, based on traditional scaling laws:

$$L(s, N) = L_0 + A \cdot S_1^{-\alpha} + B \cdot N^{-\beta} - C \cdot S_2 \cdot N^{\gamma},$$
(5)

where N is the number of non-embedding model parameters, and B, β , γ are positive constants related to N. We realize $S_2 \propto N^{\gamma}$ via a multiplier N^{γ} to the original annealing term $-C \cdot S_2$.

Fitting and Prediction with Model Size. We validate Eq. 5 by fitting the full loss curves of models with varying sizes. We then apply the obtained equation to predict full loss curve on the unseen largest model size. Results in Fig. 7b show an almost perfect fit ($R^2 > 0.998$) and prediction for entire training dynamics of larger-scale models. This indicates the effectiveness and robustness of our proposed N-extended equation. Additional N-extended experiments with other setups further confirm the robustness of our formulation (see detail in Appendix D.6).

378 4 APPLICATION

We apply our proposed formulation to validate and provide a theoretical explanation for numerous existing experimental findings regarding the training dynamics of language models. These key insights also guide researchers in selecting critical LRS before initiating model training. An interesting summary is that, *the art of learning rate schedule lies in the delicate balancing act between forward area and annealing area.*

3853864.1 DETERMINING COSINE CYCLE LENGTH AND MINIMUM LR IN COSINE LRS.

Many papers have found that in LLM pre-training 387 using cosine LRS, setting the cosine cycle length T388 as the total steps S, and setting min LR as nearly 0 389 (rather than 10% max LR) can lead to the optimal 390 loss (Hoffmann et al., 2022; Hu et al., 2024; Hägele 391 et al., 2024; Parmar et al., 2024). We theoretically 392 validate this observation using our equation in Fig. 8. The predicted loss curve with T = S and a minimum 394 LR of 0 indeed achieves the optimal loss in the final 395 step. Moreover, our equation gives a quite intuitive 396 explanation: setting T > S leads to incomplete an-397 nealing, while T < S leads to a small forward area S_1 due to early annealing. Thus, the optimal con-398 figuration is to set T equal to S. Also, setting the 399 minimum LR to 0 maximizes the annealing amount, 400 thereby increasing the annealing area S_2 , which fa-401 cilitates lower final loss. 402

(a) Learning rate curves of three types of LRS.



Figure 8: Predicted loss curves of different cycle length T and min LR in cosine LRS. The results well align with previous studies.

(b) S_1 -item and negative S_2 -item of different LRS.



4.2 IT VERIFIES AND EXPLAINS WHY WSD AND MULTI-STEP COSINE LRS ARE BETTER.

418 419 420

403

Figure 9: The comparison between S_1 -item and negative S_2 -item in different LRS.

421 Recent studies have shown that WSD LRS (Hu et al., 2024) and multi-step cosine LRS (DeepSeek-422 AI, 2024) result in lower loss compared to the traditional cosine LRS. We validate and elucidate this finding using our proposed equation. Fig. 9 shows the learning rate curve (left) and the predicted 423 loss drop (right) for different LRS. The results suggest that for WSD and multi-step cosine LRS, 424 the negative S_2 -item $(-C \cdot S_2)$ is slightly larger than that of the cosine LRS, whereas the S_1 -425 item $(A \cdot S_1^{-\alpha})$ is significantly lower. Specifically, both the WSD LRS and multi-step cosine LRS 426 unintentionally employ strategies that marginally reduces S_2 but substantially increases S_1 , leading 427 to an overall decrease in validation loss. 428

420 429 430

4.3 DETERMINING OPTIMAL ANNEALING RATIO OF WSD SCHEDULER.

In the case of WSD LRS, it is crucial to ascertain the optimal annealing ratio for training steps. Hägele et al. (2024) found that there is an optimal annealing ratio for WSD LRS, and both exces-



(a) The relationship between the predicted final loss and the ratio of annealing steps under the condition of different total steps.

(b) The relationship between predicted final loss and the forward area S_1 of different total steps. Different points denote different annealing ratios.

20k Steps 40k Steps

60k Steps 80k Steps

20.0

100k Ste

Figure 10: Illustration of the predicted loss in relation to the ratio of annealing steps and the forward area in WSD LRS (cosine annealing), presenting parabola-like curves, with a distinct optimal loss.

sively high or low annealing ratios lead to sub-optimal model performance. This phenomenon can be further elucidated through our proposed equation. Specifically, a high annealing ratio results in a significant reduction of the forward area S_1 while a low annealing ratio leads a diminished annealing area S_2 . Our scaling law equation describes the trade-off between the forward area S_1 and the annealing area S_2 about the annealing ratio.

Fig. 10 depicts the final loss predicted by our equation for various annealing ratios and total training steps. The predictions form parabola-like curves, and align well with the actual experimental results reported in previous studies. This suggests that a moderate annealing ratio, typically around 10% to 20%, is optimal, as it balances S_1 and S_2 to maximize their combined effect, thereby minimizing the overall validation loss. Moreover, our equation can directly predict the optimal annealing ratio for different total steps without large-scale experiments, which saves a lot of resources.

462 4.4 MANY OTHER TAKEAWAYS

Moreover, we use our equation to verify and explain more phenomena as follows: (1) Appendix G.1:
An empirical reason of loss dropping more sharply when LR anneals (Loshchilov & Hutter, 2016;
Ibrahim et al., 2024; DeepSeek-AI, 2024). (2) Appendix G.2: the comparison between constant and
cosine LRS, aligned with previous works (Hu et al., 2024). (3) Appendix G.3: how to choose the
optimal annealing function in WSD LRS, aligned with previous works (Hägele et al., 2024). (4)
Appendix G.4 and G.5: how to re-warmup (including re-warmup peak LR and steps) in continual
pre-training, aligned with previous works (Gupta et al., 2023). Given the instances above, we believe
that our equation can help analyze and select more training recipes in specific scenarios.

471 472

473

445

446

447

450

5 COMPARISON WITH CHINCHILLA SCALING LAW

474 5.1 REDUCTION TO CHINCHILLA SCALING LAW

Our scaling law equation can predict the full loss curve across any given LRS. In this section, we show that our equation has no contradiction with traditional scaling laws, and it is a generalized form of the Chinchilla scaling law (Hoffmann et al., 2022). Specifically, all the final loss values for different total training steps following our equation should also follow a power-law relationship. We prove this by dividing two conditions: (1) constant LRS, and (2) other LRS.

Constant LRS. In the case of a constant LRS, the annealing area S_2 is always zero and the forward area $S_1 = \eta_{max} \cdot s$, where s is the step, and η_{max} is the constant maximal LR. Thus, the whole train loss curve becomes: $L(s) = L_0 + (A \cdot \eta_{max}^{-\alpha}) \cdot s^{-\alpha} = L_0 + A' \cdot s^{-\alpha}$, which aligns with the Chinchilla scaling law equation.

485 **Other LRS.** For non-constant LRS, we use a statistical approach to show that our equation can be reduced to the Chinchilla scaling law. Specifically, we verify whether the Chinchilla scaling law ad-

507

508

521

522

Equation	LRS	Computational cost	Applicable to other LRS?
Chinchilla	Cosine	100%	No
Chinchilla	WSD (20% annealing)	21.6%	No
Chinchilla	WSD (10% annealing)	11.8%	No
Ours	Any (except constant)	$<\!\!1.0\%$	Yes

Table 2: The comparison of computational cost for fitting different scaling law equations.

equately fits the endpoints of loss curves predicted by our equation. The parameter tuple of our equation is (L_0, A, C, α) . We randomly sample different parameter tuples (detailed in Appendix E.1). Each parameter tuple represents a synthetic fitting re-

497 sult corresponding to a distinct set of experimental se-498 tups (e.g., dataset, model size, etc.). For each sampled 499 parameter tuple, we apply our equation to predict the fi-500 nal loss for different total training steps with both cosine 501 and WSD LRS, and then employ the predicted losses to 502 fit the Chinchilla scaling law. We calculate the mean and standard deviation of R^2 values for each fit. The results in 504 Table 1 demonstrate that Chinchilla scaling law fits well 505 on the data predicted by our scaling law equation. Thus,

Table 1: Mean and std of R^2 for different parameter fits.

LRS	$\mathrm{mean}(R^2)\uparrow$	$\mathbf{std}(R^2)\downarrow$
Cosine	0.972	0.056
WSD	0.979	0.053

⁵⁰⁶ our equation can be considered a generalization that can be reduced to the Chinchilla scaling law.

5.2 SCALING LAW FITTING AND PREDICTION DEMOCRATIZATION

509 Our scaling law equation allows us to utilize all loss values from a full loss curve during training, 510 while traditional scaling laws can only collect a single data point from the full loss curve. This 511 feature allows us to fit scaling laws with much less cost. For a direct comparison, we compare the 512 computational efficiency of our approach and the Chinchilla scaling law (Hoffmann et al., 2022). 513 Specifically, we assume to collect 100 data points for parameter fitting, and estimate the compu-514 tational costs needed to fit the respective scaling law equations under different LRS configurations 515 (see Table 2). More details can be found in Appendix E.2. The results indicate that our proposed 516 equation uses less than 1% of the computational cost required by the Chinchilla scaling law. Further, 517 our scaling law with LR annealing, can be universally applied to predict loss curves for unseen LRS, thus conserving even more computational resources. This approach significantly democratizes the 518 study of scaling laws in LLM pre-training, paving the way for a more environmentally friendly and 519 carbon-efficient methodology. 520

6 DISCUSSION

523 (1) We analyze the impact of the decay factor λ of our equation in Appendix H.1, and it suggests that 524 selecting a proper decay factor is important for determining the balance point between S_1 and S_2 ; 525 (2) We analyze the root reasons of the delay phenomenon mentioned in Sec. 3 in Appendix H.2. It 526 suggests that neither the Adam optimizer (Kingma & Ba, 2015) nor S_1 are the root reasons and this 527 can be an important future work; (3) We discuss some potential variation of our proposed equation 528 (e.g. $\eta = 0$ case and $L \propto S_2^{\zeta}$ variant), and investigate other possible scaling law formats with LR 529 annealing in Appendix H.3. The results validate the superiority of our proposed formula.

530 531 7 CONCLUSION

In conclusion, we propose that the loss curves of neural language models empirically adhere to a scaling law with learning rate annealing over training steps $s: L(s) = L_0 + A \cdot S_1^{-\alpha} - C \cdot S_2$. This equation can accurately predict full loss curves across unseen learning rate schedulers. We present the underlying intuition and theory for deriving our equation and demonstrate that our approach can be extended to capture the scaling effect of model sizes. Extensive experiments demonstrate that our proposed scaling law has good accuracy, scalability, and holds under various experimental setups. It also offers accurate theoretical insights to the training dynamics of LLMs, and explains numerous phenomena observed in previous studies. We believe that the scaling law with LR annealing is promising to reshape the understanding of researchers for LLM training and scaling laws.

540 REFERENCES

547

567

568

569

- Yasaman Bahri, Ethan Dyer, J. Kaplan, Jaehoon Lee, and Utkarsh Sharma. Explaining neural scaling
 laws. *Proceedings of the National Academy of Sciences of the United States of America*, 2021.
 doi: 10.1073/pnas.2311878121.
- ⁵⁴⁵ BigScience. Bloom: A 176b-parameter open-access multilingual language model. *arXiv preprint arXiv: 2211.05100*, 2022.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhari-548 wal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agar-549 wal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, 550 Daniel Ziegler, Jeffrey Wu, Clemens Winter, Chris Hesse, Mark Chen, Eric Sigler, Mateusz 551 Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec 552 Radford, Ilya Sutskever, and Dario Amodei. Language models are few-shot learners. In 553 H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin (eds.), Advances in Neu-554 ral Information Processing Systems, volume 33, pp. 1877-1901. Curran Associates, Inc., 555 2020. URL https://proceedings.neurips.cc/paper_files/paper/2020/ file/1457c0d6bfcb4967418bfb8ac142f64a-Paper.pdf. 556
- Ethan Caballero, Kshitij Gupta, I. Rish, and David Krueger. Broken neural scaling laws. *International Conference on Learning Representations*, 2022. doi: 10.48550/arXiv.2210.14891.
- Together Computer. Redpajama: an open dataset for training large language models, 2023. URL
 https://github.com/togethercomputer/RedPajama-Data.
- DeepSeek-AI. Deepseek llm: Scaling open-source language models with longtermism. arXiv preprint arXiv: 2401.02954, 2024.
- Zhengxiao Du, Aohan Zeng, Yuxiao Dong, and Jie Tang. Understanding emergent abilities of
 language models from the loss perspective. *arXiv preprint arXiv: 2403.15796*, 2024.
 - Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. The llama 3 herd of models. *arXiv preprint arXiv: 2407.21783*, 2024.
- William Fedus, Barret Zoph, and Noam Shazeer. Switch transformers: Scaling to trillion parameter
 models with simple and efficient sparsity. *arXiv preprint arXiv: 2101.03961*, 2021.
- 573 Dirk Groeneveld, Iz Beltagy, Pete Walsh, Akshita Bhagia, Rodney Kinney, Oyvind Tafjord, 574 Ananya Harsh Jha, Hamish Ivison, Ian Magnusson, Yizhong Wang, Shane Arora, David Atkin-575 son, Russell Authur, Khyathi Raghavi Chandu, Arman Cohan, Jennifer Dumas, Yanai Elazar, 576 Yuling Gu, Jack Hessel, Tushar Khot, William Merrill, Jacob Morrison, Niklas Muennighoff, 577 Aakanksha Naik, Crystal Nam, Matthew E. Peters, Valentina Pyatkin, Abhilasha Ravichander, 578 Dustin Schwenk, Saurabh Shah, Will Smith, Emma Strubell, Nishant Subramani, Mitchell Worts-579 man, Pradeep Dasigi, Nathan Lambert, Kyle Richardson, Luke Zettlemoyer, Jesse Dodge, Kyle 580 Lo, Luca Soldaini, Noah A. Smith, and Hannaneh Hajishirzi. Olmo: Accelerating the science of language models. arXiv preprint arXiv: 2402.00838, 2024. 581
- Kshitij Gupta, Benjamin Thérien, Adam Ibrahim, Mats L. Richter, Quentin Anthony, Eugene Belilovsky, Irina Rish, and Timothée Lesort. Continual pre-training of large language models: How to (re)warm your model?, 2023. URL https://arxiv.org/abs/2308.04014.
- Tom Henighan, Jared Kaplan, Mor Katz, Mark Chen, Christopher Hesse, Jacob Jackson, Heewoo
 Jun, Tom B Brown, Prafulla Dhariwal, Scott Gray, et al. Scaling laws for autoregressive generative
 modeling. *arXiv preprint arXiv:2010.14701*, 2020.
- Danny Hernandez, Jared Kaplan, Tom Henighan, and Sam McCandlish. Scaling laws for transfer.
 arXiv preprint arXiv:2102.01293, 2021.
- Joel Hestness, Sharan Narang, Newsha Ardalani, Gregory Diamos, Heewoo Jun, Hassan Kianinejad,
 Md. Mostofa Ali Patwary, Yang Yang, and Yanqi Zhou. Deep learning scaling is predictable,
 empirically. *arXiv preprint arXiv: 1712.00409*, 2017.

631

632

- Jordan Hoffmann, Sebastian Borgeaud, Arthur Mensch, Elena Buchatskaya, Trevor Cai, Eliza Rutherford, Diego de Las Casas, Lisa Anne Hendricks, Johannes Welbl, Aidan Clark, Tom Hennigan, Eric Noland, Katie Millican, George van den Driessche, Bogdan Damoc, Aurelia Guy, Simon Osindero, Karen Simonyan, Erich Elsen, Jack W. Rae, Oriol Vinyals, and Laurent Sifre. Training compute-optimal large language models. *arXiv preprint arXiv: 2203.15556*, 2022.
- Shengding Hu, Yuge Tu, Xu Han, Chaoqun He, Ganqu Cui, Xiang Long, Zhi Zheng, Yewei Fang, Yuxiang Huang, Weilin Zhao, Xinrong Zhang, Zheng Leng Thai, Kaihuo Zhang, Chongyi Wang, Yuan Yao, Chenyang Zhao, Jie Zhou, Jie Cai, Zhongwu Zhai, Ning Ding, Chao Jia, Guoyang Zeng, Dahai Li, Zhiyuan Liu, and Maosong Sun. Minicpm: Unveiling the potential of small language models with scalable training strategies. *arXiv preprint arXiv: 2404.06395*, 2024.
- 605
 Peter J. Huber. Robust Estimation of a Location Parameter. The Annals of Mathematical Statistics, 35(1):73 – 101, 1964. doi: 10.1214/aoms/1177703732. URL https://doi.org/10.1214/ aoms/1177703732.
- Alexander Hägele, Elie Bakouch, Atli Kosson, Loubna Ben Allal, Leandro Von Werra, and Martin Jaggi. Scaling laws and compute-optimal training beyond fixed training durations. *arXiv preprint arXiv: 2405.18392*, 2024.
- Adam Ibrahim, Benjamin Thérien, Kshitij Gupta, Mats L. Richter, Quentin Gregory Anthony, Eugene Belilovsky, Timothée Lesort, and Irina Rish. Simple and scalable strategies to continually pre-train large language models. *Trans. Mach. Learn. Res.*, 2024, 2024. URL https://openreview.net/forum?id=DimPeeCxKO.
- Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B. Brown, Benjamin Chess, Rewon Child,
 Scott Gray, Alec Radford, Jeffrey Wu, and Dario Amodei. Scaling laws for neural language
 models. arXiv preprint arXiv: 2001.08361, 2020.
- ⁶¹⁹ Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In Yoshua
 ⁶²⁰ Bengio and Yann LeCun (eds.), 3rd International Conference on Learning Representations, ICLR
 ⁶²¹ 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings, 2015. URL http:
 ⁶²² //arxiv.org/abs/1412.6980.
- Atli Kosson, Bettina Messmer, and Martin Jaggi. Analyzing & eliminating learning rate warmup in GPT pre-training. In *High-dimensional Learning Dynamics 2024: The Emergence of Structure* and Reasoning, 2024. URL https://openreview.net/forum?id=RveSp50ESA.
- Liyuan Liu, Haoming Jiang, Pengcheng He, Weizhu Chen, Xiaodong Liu, Jianfeng Gao, and Jiawei
 Han. On the variance of the adaptive learning rate and beyond. In 8th International Conference
 on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020. OpenRe view.net, 2020. URL https://openreview.net/forum?id=rkgz2aEKDr.
 - I. Loshchilov and F. Hutter. Sgdr: Stochastic gradient descent with warm restarts. *International Conference on Learning Representations*, 2016.
- 634 Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. *arXiv preprint arXiv:* 635 *1711.05101*, 2017.
- Eric J. Michaud, Ziming Liu, Uzay Girit, and Max Tegmark. The quantization model of neural scaling. In Alice Oh, Tristan Naumann, Amir Globerson, Kate Saenko, Moritz Hardt, and Sergey Levine (eds.), Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 16, 2023, 2023. URL http://papers.nips.cc/paper_files/paper/2023/hash/5b6346a05a537d4cdb2f50323452a9fe-Abstract-Conference.html.
- Niklas Muennighoff, Alexander M Rush, Boaz Barak, Teven Le Scao, Nouamane Tazi, Aleksandra Piktus, Sampo Pyysalo, Thomas Wolf, and Colin Raffel. Scaling data-constrained language models. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023. URL https://openreview.net/forum?id=j5BuTrEj35.
- Jorge Nocedal. Updating quasi newton matrices with limited storage. *Mathematics of Computation*, 35(151):951–958, July 1980. ISSN 0025-5718. doi: 10.1090/S0025-5718-1980-0572855-7.
 - 12

648 649	OpenAI. Gpt-4 technical report. arXiv preprint arXiv:2303.08774, 2023.
650	Juninder Parmar, Saniev Satheesh, Mostofa Patwary, Mohammad Shoevhi, and Bryan Catanzaro
651	Reuse, don't retrain: A recipe for continued pretraining of language models. <i>arXiv preprint</i>
652	arXiv: 2407.07263, 2024.
653	
654 655	Raffel, Leandro Von Werra, and Thomas Wolf. The fineweb datasets: Decanting the web for the finest text data at scale. <i>arXiv preprint arXiv: 2406.17557</i> , 2024.
656 657 658	Leslie N Smith. Cyclical learning rates for training neural networks. In 2017 IEEE winter conference on applications of computer vision (WACV), pp. 464–472. IEEE, 2017.
659	Samuel I. Smith Diatar Jan Kindermans, Chris Ving, and Quee V. Le. Don't decay the learning rate
660	increase the batch size. In 6th International Conference on Learning Representations. ICLR 2018
661 662	Vancouver, BC, Canada, April 30 - May 3, 2018, Conference Track Proceedings. OpenReview.net, 2018. URL https://openreview.net/forum?id=B1Yy1BxCZ.
663 664 665 666	Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Niko- lay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open founda- tion and fine-tuned chat models. <i>arXiv preprint arXiv:2307.09288</i> , 2023.
667 668	Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. <i>arXiv preprint arXiv: 1706.03762</i> , 2017
669	2017.
670	
670	
673	
674	
675	
676	
677	
678	
679	
680	
681	
682	
683	
684	
685	
686	
687	
688	
689	
690	
602	
692	
694	
695	
696	
697	
698	
699	
700	
701	

702 A IMPACT OF WARMUP STEPS

We conduct experiments on the impact of learning rate warmup steps. As shown in Fig. 11, we find that 500 warmup steps can speed up convergence, and get the lowest validation loss compared to 100 or no LR warmup. The finding is aligned with previous works Liu et al. (2020); Kosson et al. (2024). The experimental results also guide us to choose 500 warmup steps in the main experiments of this work ².



Figure 11: The comparison of the true loss curve of different warmup steps. We experiment on cosine LRS with 20K total steps.

B FITTING DETAILS

Given a learning rate scheduler, we can easily compute out S_1 and S_2 of each step in advance. To estimate (L_0, A, C, α) , we adopt a similar fitting method as Chinchilla scaling law (Hoffmann et al., 2022). Specifically, we minimize the Huber loss (Huber, 1964) between the predicted and the observed log loss using the L-BFGS algorithm (Nocedal, 1980):

$$\min_{L_0, A, C, \alpha} \quad \sum_{\text{Step } i} \text{Huber}_{\delta} \left(\log \hat{L}(i) - \log L(i) \right)$$
(6)

We implement this by the utilization of minimize in scipy library. Huber loss is to enhance to robustness of the fitting results and we set δ of Huber loss as 1.0×10^{-3} . We mitigate the potential issue of local minima of fitting by choosing the optimal fit from a range of initial conditions. Note that in practice, we can also fit the full loss curves using multiple LRS with a single tuple of (L_0, A, C, α) . In this situation, we sum the Huber losses in Eq. 6 of all fitted LRS.

C EXPERIMENTAL SETUPS

In this work, we use multiple sets of experimental setups, in order to validate that our equation can work across different experimental setups. For clarification, we present the experimental setup list as shown in Table 3. In our most experiments, we use the main setting A. Other than the five settings, we also successfully fit our equation on BLOOM's and OLMo's loss curves, and their experimental settings are totally different. Refer to their papers for the experimental settings (BigScience, 2022; Groeneveld et al., 2024).

 ²Note LR warmup in training from scratch is different from LR re-warmup in continual training, where we
 do not regard re-warmup steps as a hyper-parameter and will show how to apply our equation to find optimal re-warmup recipes in Appenidx. G.4 and G.5.

Table 3: Experimental settings adopted in this work. Model size denotes the number of non-embedding paramters. Our datasets include Fineweb (Penedo et al., 2024) and RedPajama-CC (Computer, 2023). * denotes pre-training multilingual dataset including mixture of sources such as common crawls, books, arxiv, code, etc. We use AdamW Optimizer (Kingma & Ba, 2015; Loshchilov & Hutter, 2017), denoted as AO. Most experiments adopt Llama-3's tokenizer (Dubey et al., 2024). Ext Llama-2's is extended from Llama-2's tokenizer (Touvron et al., 2023) by adding vocabulary.

773	Setups	Setting A (mainly)	Setting B	Setting C
774	Model Size	594 M	293M	multiple
775	Train Dataset	Fineweb	Finweb	Mixture-train*
776	Val Dataset	RedPajama-CC	RedPajama-CC	Mixture-valid*
777	Total Steps	60K	120 K	143 K
778	Maximal LR	2×10^{-4}	2×10^{-4}	$1.381 imes 10^{-3}$
779	Warmup Steps	500	100	500
780	Batch Size (tokens)	4M	2 M	4 M
781	Sequence Length	4096	4096	4096
782	Tokenizer	Llama-3's	Llama-3's	Ext Llama-2's
783	eta_1,eta_2 in AO	0.9, 0.95	0.9, 0.95	0.9, 0.95
784	Weight Decay	0.1	0.1	0.1
785	Gradient Clip	1.0	1.0	1.0
786	Setups	Setting D (MoE)	Setting E (1.4T tokens)	
	A	0	8	
787	Model Size	$8 \times 106 M$	1704M	
787 788	Model Size Train Dataset	$8 \times 106M$ Fineweb	1704M Mixture-train*	
787 788 789	Model Size Train Dataset Val Dataset	8 × 106M Fineweb RedPajama-CC	1704M Mixture-train* Mixture-valid*	
787 788 789 790	Model Size Train Dataset Val Dataset Total Steps	8 × 106M Fineweb RedPajama-CC 60K	1704M Mixture-train* Mixture-valid* 350K	
787 788 789 790 791	Model Size Train Dataset Val Dataset Total Steps Maximal LR	$8 \times 106M$ Fineweb RedPajama-CC 60K 2×10^{-4}	$1704M$ Mixture-train* Mixture-valid* $350K$ 6×10^{-4}	
787 788 789 790 791 792	Model Size Train Dataset Val Dataset Total Steps Maximal LR Warmup Steps	$8 \times 106M$ Fineweb RedPajama-CC 60K 2×10^{-4} 500	$1704M$ Mixture-train* Mixture-valid* $350K$ 6×10^{-4} 1000	
787 788 789 790 791 792 793	Model Size Train Dataset Val Dataset Total Steps Maximal LR Warmup Steps Batch Size (tokens)	$8 \times 106M$ Fineweb RedPajama-CC $60K$ 2×10^{-4} 500 $4M$	$1704M$ Mixture-train* Mixture-valid* $350K$ 6×10^{-4} 1000 $4M$	
787 788 789 790 791 792 793 794	Model Size Train Dataset Val Dataset Total Steps Maximal LR Warmup Steps Batch Size (tokens) Sequence Length	$8 \times 106M$ Fineweb RedPajama-CC 60K 2×10^{-4} 500 4M 4096	$1704M$ Mixture-train* Mixture-valid* $350K$ 6×10^{-4} 1000 $4M$ 8192	
787 788 789 790 791 792 793 794 795	Model Size Train Dataset Val Dataset Total Steps Maximal LR Warmup Steps Batch Size (tokens) Sequence Length Tokenizer	$8 \times 106M$ Fineweb RedPajama-CC $60K$ 2×10^{-4} 500 $4M$ 4096 Llama-3's	$1704M$ Mixture-train* Mixture-valid* $350K$ 6×10^{-4} 1000 $4M$ 8192 Llama-3's	
787 788 789 790 791 792 793 794 795 796	Model SizeTrain DatasetVal DatasetTotal StepsMaximal LRWarmup StepsBatch Size (tokens)Sequence LengthTokenizer β_1, β_2 in AO	$\begin{array}{c} 8 \times 106 \mathrm{M} \\ \mathrm{Fineweb} \\ \mathrm{RedPajama-CC} \\ 60 \mathrm{K} \\ 2 \times 10^{-4} \\ 500 \\ 4 \mathrm{M} \\ 4096 \\ \mathrm{Llama-3's} \\ 0.9, 0.95 \end{array}$	$1704M$ Mixture-train* Mixture-valid* 350K 6×10^{-4} 1000 $4M$ 8192 Llama-3's 0.9, 0.95	
787 788 789 790 791 792 793 793 794 795 796 797	Model SizeTrain DatasetVal DatasetTotal StepsMaximal LRWarmup StepsBatch Size (tokens)Sequence LengthTokenizer β_1, β_2 in AOTop-k Experts	$8 \times 106M$ Fineweb RedPajama-CC 60K 2×10^{-4} 500 4M 4096 Llama-3's 0.9, 0.95 2	1704M Mixture-train* Mixture-valid* 350K 6×10^{-4} 1000 4M 8192 Llama-3's 0.9, 0.95	
787 788 789 790 791 792 793 793 794 795 796 797 798	Model SizeTrain DatasetVal DatasetTotal StepsMaximal LRWarmup StepsBatch Size (tokens)Sequence LengthTokenizer β_1, β_2 in AOTop-k ExpertsAuxiliary Loss	$8 \times 106M$ Fineweb RedPajama-CC 60K 2×10^{-4} 500 4M 4096 Llama-3's 0.9, 0.95 2 0.01	1704M Mixture-train* Mixture-valid* 350K 6×10^{-4} 1000 4M 8192 Llama-3's 0.9, 0.95 -	
787 788 789 790 791 792 793 794 795 796 797 798 799	Model Size Train Dataset Val Dataset Total Steps Maximal LR Warmup Steps Batch Size (tokens) Sequence Length Tokenizer β_1,β_2 in AO Top-k Experts Auxiliary Loss Weight Decay	$8 \times 106M$ Fineweb RedPajama-CC 60K 2×10^{-4} 500 4M 4096 Llama-3's 0.9, 0.95 2 0.01 0.1	1704M Mixture-train* Mixture-valid* 350K 6×10^{-4} 1000 4M 8192 Llama-3's 0.9, 0.95 - - 0.1	

B10 D OUR SCALING LAW ON EXTENSIVE EXPERIMENTS SETUPS B11

D.1 ANOTHER SET OF TRAINING HYPER-PARAMETERS

Fig. 2 and Fig. 3 show that our equation can work very well on our main experimental setup. For proving that our scaling law with LR annealing can apply to different (but given) experimental settings, we change the setting from A to B (refer to Table 3) and observe whether our equation can still work or not. The fitting results are shown in Fig. 12. The prediction results are shown in Fig. 13. The results suggest that our scaling law with LR annealing can still work well across different experimental setups.

819 820 821

822

812

813 814

815

816

817

818

D.2 EXPERIMENTS ON ANOTHER ARCHITECTURE: MOE

823 Fig. 2 and Fig. 3 show that our equation can work very well on the dense Llama-like architec-824 ture (Vaswani et al., 2017; Touvron et al., 2023). We prove that our scaling law can also apply to 825 different model architectures and we replace Dense model with Mixture of Experts (MoE) architec-826 ture. We add widely-used auxiliary loss to do load balancing among experts (Fedus et al., 2021). 827 The experimental setting is shown as Setting D in Table 3. Moreover, we change the LRS and total 828 steps to 60K WSD with 10K annealing steps in fitting, testing whether our scaling law is effective 829 under various circumstances. The fitting results are shown in Fig. 14 while the prediction results are 830 shown in Fig. 15. The results suggest that our scaling law can with LR annealing can still work well on MoE architecture. 831

832 833

834

D.3 SCALING UP: PREDICTION FOR MUCH LONGER STEPS

Our equation has proven its utility in predicting the validation loss over a significantly large number of total steps. This scalability feature is particularly useful in handling large-scale training scenarios.

To illustrate its effectiveness, we apply our equation to predict the loss curve during the annealing stage of the training process. The model we train is a sizable 1.7 billion parameter model, and the training involved a tremendous number of 1,400 billion total training tokens. This is a considerable scale that tests the practicality and effectiveness of our equation. The specific experimental setup is Setting E, which can be found in Table 3.

The fitting and prediction results are shown in Figure 16 and Figure 17 respectively. It shows that
we successfully get to know the loss curve in the critical annealing stages after 10x longer steps in
advance, which is crucial to handle the relationship between training dynamics and training recipes.
For example, Llama-3 adopts annealing to do pre-training data selection (Dubey et al., 2024).

847 848

849

D.4 OPEN-SOURCED FULL LOSS CURVES

For further verification for our proposed scaling law, we apply our equation on open-sourced language models and the corresponding full loss curves, including BLOOM-176B (BigScience, 2022)
and OLMo-1B (Groeneveld et al., 2024). As shown in Fig. 18, our equation also fits very well on
the open-sourced model training curves, even when the model size scales up to 176B (e.g. BLOOM)
and token number scales up to 2000B over 740K steps (e.g. OLMo).

855 856

857

D.5 Ablation Studies on S_1 and S_2

In Sec. 3, we present the strong capability of our proposed scaling law. The formulation of our scaling, $L(s) = L_0 + A \cdot S_1^{-\alpha} - C \cdot S_2$, contains two key components including S_1 and S_2 . In this section, we conduct ablation studies on S_1 and S_2 . Specifically, we compare the forms without S_1 or S_2 using setting A. For each format, we re-fit the full loss curves under 20K cosine + constant and re-predict the full curves on longer steps under different LRS. The prediction error results are shown as Table 4. The results indicate that the prediction error increases significantly in the absence of either S_1 or S_2 and suggest that each component in our scaling law is important and indispensable.



Figure 12: Full loss curve **fitting** on cosine (30K steps to $\eta_{min} = 0$) and constant LRS. The figures omit the warmup in the first 100 steps. After fitting, we get a **universal** loss equation $L = 2.761 + 0.517 \cdot S_1^{-0.491} - 0.458 \cdot S_2$. Refer to setting B in Table 3 for experimental setups.



(c) Full curve prediction of WSD LRS (110K total steps; 10 % cosine annealing to $\eta_{min} = 0$).

Figure 13: Full loss curve **prediction** (120K steps) by the universal loss curve equation across various LRS, fitted in Fig. 12. The left, the medium, and the right figures in each row are learning rate curve, zoomed-out loss prediction, and zoomed-in loss prediction, respectively. The red rectangle means the zoomed-in zone. The figures omit the warmup in the first 100 steps. Please note that these are predictive results, which means that none of the points in this figure (except constant LRS) are involved in the fitting process. The mean prediction errors across various LRS are low to $\sim 0.2\%$. Refer to setting B in Table 3 for experimental setups.



Figure 15: Full loss curve **prediction** on MoE model by the universal loss curve equation across various unseen LRS fitted in Fig. 14. The left, the medium, and the right figures in each row are learning rate curve, zoomed-out loss prediction, zoomed-in loss prediction, respectively. The red rectangle means the zoomed-in zone. The LR curve figures omit 500 warmup steps. Note that these are all predictive results, and none of the points in the figures are involved in the fitting process. The mean prediction errors across various LRS are low to $\sim 0.2\%$. Refer to setting D in Table 3 for experimental setups.



Figure 16: Full loss curve **fitting** on 30K Steps. After fitting, we get a **universal** loss equation $L = 2.788 + 0.906 \cdot S_1^{-0.416} - 0.254 \cdot S_2$. Refer to setting E in Table 3 for experimental setups.



Figure 17: Full loss curve **prediction** (350K steps) by the universal loss curve equation under WSD LRS (10% cosine annealing ratio to $\eta_{min} = 0$). We adopt our equation and accurately predict the loss curve in the annealing stage after the 10x longer steps. This is meaningful to the development for large-scale LLM pre-training. Refer to setting E in Table 3 for experimental setups.



(a) Full loss curve fitting on BLOOM-176B.

(b) Full loss curve fitting on OLMo-1B (2T tokens).

Figure 18: Open-sourced full loss curve fitting using our proposed equation, which shows that our equation has strong scalability on model size and token number. We extract the curve of BLOOM from https://huggingface.co/bigscience/bloom/tensorboard, and we choose the column lm-loss-validation/valid/lm loss validation as validation loss. We extract the curve of OLMo from https://wandb.ai/ai2-llm/OLMo-1B?
 nw=nwuserdirkgr, and we choose the column eval/pile/CrossEntropyLoss as validation loss. Both models adopt cosine LRS.

	Scaling Law Forms	Cosine	Multi-step Cosine	WSD	Cyclic
	$L(s) = L_0 + A \cdot S_1^{-\alpha} - C \cdot S_2$ w/o S ₁ : $L(s) = L_0 + A \cdot s^{-\alpha} - C \cdot S_2$ w/o S ₂ : $L(s) = L_0 + A \cdot S_1^{-\alpha}$	0.159% 0.465% 1.261%	0.176% 0.458% 1.265%	0.235% 1.162% 1.519%	0.322% 1.139% 1.203%
		1120170	1.20070	1.01770	1120070
	3.9 -		× Ground-truth	LOSS	
,	3.8 - 3.7 - Fitting Curve: Loss = 1.41	2 + 0.921*S1^(-0.7	761) N=18.34M		
	3.6 + 2.562*N^(-0.137) - 0.0	69*S2*N^(0.117)	N=106 2M		
	$R^2 = 0.99837$		N=161.9M		
	3.4 -		N=293M		
	ss 3.3		N=401M		
			N=594M		
			N=1720M		
		000000000000000000000000000000000000000	Norman .		
	> 2.8 -			Miciacepee	
	2.7				
	2.6				

	2.3	20000000000000000000000000000000000000		xxxxxxxxxxx	
	2.2			*******	
	0 20000 40000	60000 80	000 100000 120000	140000	
		Steps			

Table 4: Ablation studies on S_1 and S_2 in our scaling law formulation. The prediction errors on different LRS for each form are reported.

Figure 19: Curve fitting on cosine LRS (143K steps to $\eta_{min} = 0$) of many model sizes using our scaling law extended to model size N. Refer to setting C in Table 3 for experimental setups.

1055 1056

1052

1057 D.6 OUR *N*-EXTENDED SCALING LAW ON ANOTHER EXPERIMENTS SETUPS

Fig. 7b show that our N-extended equation can work very well on our main experimental setup. Similarly, for proving that our N-extended scaling law can apply to different (but given) experimental settings, we change the setting from A to C (refer to Table 3) and observe whether our equation can still work or not. The fitting results are shown in Fig. 19. The results suggest that our N-extended scaling law with LR annealing can still work well across different experimental setups.

1063

1064

E COMPARISON WITH CHINCHILLA SCALING LAW

E.1 REDUCTION TO CHINCHILLA SCALING LAW

1068 We have proved that our scaling law can be reduced to chinchilla scaling law for constant LRS in 1069 Sec. 5. For other learning rate schedulers, we adopt a method based on statistics to show that our 1070 scaling law function can be reduced to the chinchilla scaling law. Specifically, we check whether chinchilla scaling law fits well the endpoints of loss curves predicted by our scaling law. The pa-1071 rameter tuple of our equation is (L_0, A, C, α) . We then randomly sample 1000 sets of param-1072 eter tuples in some uniform distributions: $L_0 \sim U(1,3), A \sim U(0.3, 0.5), C \sim U(0.2, 0.6),$ 1073 $\alpha \sim U(-0.6, -0.4)$. Each parameter tuple could be seen as the fitting result of a distinct set of ex-1074 perimental setups ³ (e.g. dataset, batch size, model size, etc.). For each generated parameter tuple, 1075 we apply our equation to predict the final loss of different total training steps on two LRS including 1076 cosine and WSD (10% annealing ratio). range from 5K steps to 60K steps. We conduct the predic-1077 tion on two LRS including cosine and WSD (10% annealing ratio). The predicted final loss points

³It's worth noting that some of these sampled parameter tuples might not be reasonable or likely to happen in real-world scenarios, but we choose to keep them nonetheless.



(a) The predicted loss of different total steps with cosine LRS and the fitted chinchilla curve.



Figure 20: Chinchilla scaling law fits well the validation loss endpoints predicted by our formulation,
 taking cosine LRS (on the left) and WSD LRS (on the right) as examples.

are used to fit the chinchilla equation through minimizing the Huber loss. The fitting examples areshown in Fig. 20.

1102 E.2 SCALING LAW COMPUTATIONAL COST COMPARISON

1104 We suppose a scenario where it requires 100 fitting points to get the parameters of scaling laws. We 1105 assume the distance between each point as K steps. We compute the required training steps using 1106 different approaches as follows:

- Adopting Chinchilla scaling law, typical cosine LRS requires total steps of at least $1K + 2K + 3K + \cdots + 100K = 5050K$;
- Adopting Chinchilla scaling law, WSD LRS (notating annealing ratio as r) requires total steps of at least $(1K + 2K + 3K + \dots + 100K)r + 100K(1 r) = (100 + 4950r)K$.
- Adopting our scaling law, all we need is only one training curve with moderate total steps (and the number of fitting points is far more than 100), such as one curve with 50K steps ⁴

F WSD SCHEDULER AND ANNEALING FUNCTIONS

Hu et al. (2024) proposes a warmup-stable-decay (WSD) LRS including three learning rate stages, which could help get a lower validation loss compared to the typical cosine LRS. The format is like

1123

1129

1130 1131

1093

1094

1098

1101

1103

1107

1108

1109

1110

1111 1112

1113

1114 1115

1116 1117

1118

$$WSD(s) = \begin{cases} \frac{s}{T_{\text{warmup}}} \eta_{max}, & s \le T_{\text{warmup}} \\ \eta_{max}, & T_{\text{warmup}} < s \le T_{\text{stable}} \\ \eta_{min} + f(s) \cdot (\eta_{max} - \eta_{min}), & T_{\text{stable}} < s \le T_{\text{total}} \end{cases}$$
(7)

1124 1125 Where $0 \le f(s) \le 1$ is typically a decreasing function about step s, and η_{max} is the maximal learn-1126 ing rate. Hägele et al. (2024) consolidates the effectiveness of WSD scheduler by many empirical 1127 experiments. Moreover, Hägele et al. (2024) also finds that using 1-sqrt annealing and a moderate 1128 annealing ratio (e.g. 20%) can further decrease the final loss. The 1-sqrt annealing is defined as:

$$f(s) = 1 - \sqrt{\frac{s - T_{\text{stable}}}{T_{\text{total}} - T_{\text{stable}}}}$$
(8)

⁴The empirical rule that more fitting points always achieve better fitting results always holds true. Our equation can also use more points and LRS for fitting, such as 30K constant + 70K cosine. Nevertheless, we can collect far more fitting points than the typical scaling law with significantly fewer training steps.





Figure 22: How S1-item and negative S_2 -item changes in a WSD scheduler. Gray area means the amount of loss drop brought by S_1 and S_2 in annealing stage.



Figure 23: Comparison of constant and cosine LRS in different steps.

early stages, the neural network is exploring globally and it is a suitable time to use a larger LR; In the later stages, the neural network is exploring locally and it is a suitable time to use a smaller LR.

G.3 IT VERIFIES AND EXPLAINS THAT THE OPTIMAL ANNEALING FUNCTION IN WSD LRS DEPENDS ON THE ANNEALING RATIO.

In the context of the WSD LRS, the selection of the annealing method in the annealing stage is also pivotal to optimize the training process. Hägele et al. (2024) conclude that the 1-sqrt annealing (refer to Appendix F for 1-sqrt function and curve) yields a lower final loss compared to the other annealing methods (e.g. cosine). They claim that the conclusion holds true across different annealing ratios.

However, as we predict using our equation (Fig. 24a), the 1-sqrt annealing approach does get a lower loss than the cosine annealing approach when using small annealing ratios (e.g. 10%), but it performs much worse than the cosine annealing approach when using 50% annealing ratio.

To verify whether the predictions from our equation are accurate, we conduct experiments by training models using different annealing methods and ratios within a fixed 50K total steps. As illustrated in Fig. 24b, at a 10% annealing ratio, the 1-sqrt method outperforms the cosine method, whereas at a 50% annealing ratio, the latter method exhibits a lower final loss. The true experimental results align quite well with our prediction, which also overturns some of the conclusions made by previous works. We conclude that the optimal annealing function in WSD LRS depends on the annealing ratio.

1241 Our scaling law function provides an explanatory framework for these observations. We draw the LR curves of 1-sqrt and cosine annealing in Appendix F. At 10% annealing ratio, although the forward

(a) The predicted loss curve of cosine and 1-sqrt annealing method of different annealing ratio.

(b) The **true** loss curve of different annealing ratios with cosine and 1-sqrt annealing methods.

Figure 24: The predicted (left) and true loss (right) of cosine and 1-sqrt annealing method at different annealing ratios. Experimental results (right), aligned with our prediction (left), refute the previous finding "the order and results of different annealing hold across settings" (Hägele et al., 2024).

(a) The predicted validation loss with different rewarmup max LR in the continual pre-training process. All the re-warmup steps are 500 steps.

(b) The predicted validation loss with different rewarmup max learning rate in the continual pretraining process.

Figure 25: The predicted validation loss with different re-warmup max learning rate and re-warmup steps in the continual pre-training process. The LRS of continual pre-training is cosine (T=100K) and the min learning rate is 0.

1280 area S_1 of the cosine method is slightly larger than that of the 1-sqrt method, the larger annealing 1281 area S_2 of the 1-sqrt method plays a more critical role in reducing the overall final loss. However, 1282 as the annealing ratio increases, the difference of S_1 between two LRS gradually becomes larger 1283 and larger, till breaking the delicate balance between S_1 and S_2 at 50% annealing ratio, resulting in 1284 a lower final loss for the cosine method. This relationship underscores the importance of carefully selecting the annealing strategy to optimize model training outcomes within the WSD scheduler. 1285 Still, our equation can help predict a better annealing method without experiments, which saves a 1286 lot of resources. 1287

1288

1279

1255

1259

G.4 IT VERIFIES AND EXPLAINS THAT IN CONTINUAL PRE-TRAINING, THE HIGHER MAX LEARNING RATE TO RE-WARMUP, THE HIGHER THE INITIAL PEAK LOSS WILL BE, AND THEN THE MORE SHARPLY IT WILL DECREASE.

1292

In continual pre-training (CPT), the learning rate scheduler is usually set as re-warmup to a new max LR at the beginning. By many experiments, Gupta et al. (2023) concludes that the higher max learning rate to re-warmup, the higher the initial peak loss will be, and then the more sharply it will decrease.

1296 According to our scaling law function ⁵, in the re-warmup process, the annealing area S_2 will reduce 1297 to a negative value ($S_2 < 0$) and thus the validation loss increases. The higher max LR in re-warmup, 1298 the annealing area S_2 becomes more negative and thus there would be a higher peak loss. But still, 1299 higher max LR could make the forward area S_1 grow faster and the loss decreases more sharply after 1300 re-warmup. We use the fitted equation to predict the continual pre-training process with different 1301 max LR as shown in Fig. 25a. The predicted loss curves reproduce a quite similar phenomenon with 1302 previous works (Gupta et al., 2023).

There is a more profound strategy using our equation in CPT. As shown in Fig. 25a, after ensuring total steps during CPT, we can apply our equation to predict a better max LR and scheduler to get the lowest final loss without experiments, which saves a lot of resources.

- 1306
- 1307 1308 1309

G.5 IT VERIFIES AND EXPLAINS THAT IN CONTINUAL PRE-TRAINING, THE STEPS OF RE-WARMUP HAVE LITTLE IMPACT ON THE FINAL LOSS.

Meanwhile, how many steps to re-warmup is another important issue in the continual pre-training. Gupta et al. (2023) find that the longer re-warmup steps could smooth the transition of loss curve but the number of re-warmup steps does not significantly influence the final validation loss. We use the fitted equation to predicted the continual pre-training dynamics with different re-warmup steps. The results, shown in Fig. 25b, present a good alignment with previous works (Gupta et al., 2023).

Based on our theory, given the fixed max LR, when the re-warmup steps are longer, the annealing area decreases more slowly and the loss curve rises more smoothly, but both final S_1 and S_2 are quite stable across different re-warmup steps. First, the annealing area S_2 of different re-warmup steps are very close due to the same max LR and the same min LR. Besides, though different re-warmup steps bring in temporary distinct losses, re-warmup only cover a small percentage compared with all training steps. Thus, the forward area S_1 is also close across different re-warmup steps, resulting in the close overall loss across different steps of re-warmup.

1321 1322 1323

1324

1325

H DISCUSSION

H.1 The impact of Decay Factor λ

¹³³⁶ 1337

Figure 26: The comparison of fitting effect of different decay factor λ .

1338 The decay factor λ in our equation plays a crucial role to indicate the information retaining degree 1339 in LR annealing. We set λ as 0.999 in our all experiments. We explore the difference from another 1340 decay factor $\lambda = 0.99$. We fit and get different equations for different λ . We compare them (1) on 1341 the predicted loss curves for 1-square and 1-sqrt annealing methods, and (2) on the predicted loss 1342 curves in different annealing ratios of WSD LRS (cosine annealing).

The results, illustrated in Fig. 26 and 27, reveal several key insights into the impact of decay factor:

Delay Steps. A larger decay factor results in longer delay steps. Comparing Fig. 26b and Fig. 26c, $\lambda = 0.999$ introduces a more obvious delay phenomenon, which is consistent across both the 1-

⁵Strictly speaking, continual pre-training process often include LR re-warmup as well as data distribution shift. Here we primarily research on the condition where there is no distribution shift between two training stages. The conclusions transfer across most cases because the loss change brought by LR re-warmup is significantly larger than the loss change brought by data distribution shift (Gupta et al., 2023; Ibrahim et al., 2024).

Figure 27: The predicted loss in different annealing ratios of WSD LRS for $\lambda = 0.99$ and $\lambda = 0.999$.

1368 1369 square and 1-sqrt annealing methods. The root reason is simple: larger λ can retain more LR historical momentum, causing longer delay steps after LR finish annealing.

1371 **Optimal Annealing Ratio.** a larger decay factor tends to favor a higher annealing ratio. As shown 1372 in Fig. 27, The optimal annealing ratio of $\lambda = 0.999$ is larger than that of $\lambda = 0.99$. Meanwhile, 1373 due to the similar reason, $\lambda = 0.999$ favors 1-sqrt annealing method while $\lambda = 0.99$ favors 1-square 1374 annealing method, as shown in Fig. 26.

Balance Point between S_1 and S_2 . More essentially, the selection of λ decides the balance point of S_1 and S_2 . For example, $\lambda = 0.999$ means that, LR annealing only retain the information of previous approximately $\frac{1}{1-\lambda} = 1000$ steps, which can be seen as the window size of LR annealing momentum. The window size could be very close to the optimal annealing steps. After reaching window size, S_2 increases very slowly, with the cost of large decrease of S_1 .

¹³⁸⁰ The analyses above highlights the importance of selecting a decay factor that aligns closely with ¹³⁸¹ empirical data to ensure the accuracy of predictions. We recommend that the future developers try ¹³⁸² different λ for their own setups ⁶.

1383

H.2 Possible Root Reasons of Delay Phenomenon in Learning Rate Annealing
 1385

In Sec. 3, we discover the delay phenomenon, which proves that LR annealing has momentum. Wediscuss possible root reasons of the phenomenon in this section.

Adam Optimizer? No. We notice that Adam optimizer (Kingma & Ba, 2015) also has the firstorder momentum decay factor β_1 and the second-order momentum decay factor β_2 , which presents the possible connection to the the delay phenomenon.

1391 1392 We keep $\beta_1 = 0.9$, and conduct delay experiments on different $\beta_2 \in \{0.95, 0.99, 0.999\}$ (default: 1393 0.95) of AdamW optimizer (Loshchilov & Hutter, 2017) to observe whether larger β_2 causes a 1394 more longer delay steps. The learning rate and ground-true loss curve are shown in Fig. 28a. It 1395 suggests that the ground-truth loss curves of different β_2 almost coincide with each other, and their 1396 delay steps are also the same. Therefore, we believe that Adam optimizer has little to do with the 1397 delay phenomenon, despite its momentum form seeming very related to our experiments. Speaking 1398 momentum, attempting to discover a connection between them, but the fitting results were a mess.

Forward Area S_1 ? Not Really. No matter how LR changes, S_1 is always increasing over steps, resulting in consistently reducing the validation loss brought from S_1 . Therefore, the forward area,

^{1402 &}lt;sup>6</sup>Actually, λ can be fitted as a parameter, instead of a hyper-parameter requiring manual tuning. We regard 1403 λ as a hyper-parameter because $\lambda = 0.999$ performs well in our all experiments. Besides, fitting with λ could bring in additional time complexity due to the recomputation of S_2 .

(a) The comparison of true loss curve with setting dif-1416 ferent β_2 of Adam optimizer. 1417

2.0

1.5

ι.0

0.0

(10)

Figure 28: The possible root reason analysis (left: Adam optimizer, right: S_1) of delay phenomenon.

1421 S_1 would lengthen delay steps in LR annealing then re-warmup, but would shorten delay steps in 1422 LR re-warmup then annealing. The delay phenomenon is indeed related to S_1 . 1423

But still, S_1 is not all the reasons of delay phenomenon. We prove this by Fig. 5b, which suggests that 1424 even though in LR re-warmup then annealing, the delay phenomenon, while not that pronounced, 1425 still exists. Moreover, we conduct delay experiments by adjusting the slope of LR after tuning point 1426 of LR. As shown in Fig. 28b, We find that more smooth slope of LR re-warmup, with smaller S_1 , 1427 but still causes longer delay steps. Therefore, we conclude that S_1 indeed influences the specific 1428 delay length, but is not the root reason. 1429

Other Possible Reasons? The delay phenomenon could be intuitive in some cases. For example, 1430 suppose that learning rate decreases directly from 2e-4 to 2e-5 in one step, and maintains 2e-5. In 1431 this case, although the loss would decrease to a lower value but the parameter changes in one step 1432 is too small in neural networks. Given a sudden low LR, neural networks still require some steps 1433 to gradually optimize to a local minimum, incurring delay phenomenon. But still, analysis above 1434 still ends with a rough description, and we have not figured out the root reasons and look forward to 1435 completing this part in future work.

1436 1437

1418

1419 1420

1438 1439

OTHER POSSIBLE SCALING LAW FORMATS WITH LR ANNEALING H.3

Adding a LR-weighted Coefficient to S_2 ? Imagine that when LR anneals to nearly 0, the neural 1440 network's parameters almost do not change and the validation loss should not change, either. How-1441 ever, as defined in our equation, Eq. 1, S_2 still has historical momentum even if LR is nearly 0, 1442 making the loss continue to decrease and misalign with observed training dynamics. 1443

To cover this corner case, we try a revision to our equation and add a LR-weighted coefficient to 1444 S_2 . Specifically, we adjust S_2 to more approach 0 when η is close to 0, counteracting the original 1445 formulation's tendency to overestimate loss reduction when $\eta \approx 0$. 1446

1 (m

1447 The revised equation for the annealing area S_2 in our scaling law function is as follows: 1448

1449

1449
1450
1451
1452

$$m_i = \lambda \cdot m_{i-1} + (\eta_{k-1} - \eta_k)$$

 $= \sum_{k=1}^i (\eta_{k-1} - \eta_k) \cdot \lambda^{i-k}$

1453

1454

- $S_2 = \sum_{i=1}^{s} m_i \cdot \boldsymbol{\eta}_i^{\boldsymbol{\epsilon}}$ 1455
- 1456

Where the red part is the added LR-weighted coefficient and ϵ is a undetermined positive constant. 1457 ϵ could be very small in practice.

We have tried the revised function to fit data. We find that the fitting results are quite similar and ϵ is very close to 0, showing little use in practical effect. Hence, we adopt the original format in our experiments ⁷.

1475 $L \propto S_2^{\zeta}$ rather than $L \propto S_2$? Actually, all we know is that L and S_2 have a positive correlation. 1476 Thus $L \propto S_2^{\zeta}$ rather than $L \propto S_2$ might be a more reasonable assumption. That is, our equation 1477 would be changed to $L(s) = L_0 + A \cdot S_1^{-\alpha} - C \cdot S_2^{\zeta}$. Theoretically, the introduction of ζ as an 1478 additional fitting parameter is expected to provide a more nuanced control over how changes in the 1479 learning rate annealing affect validation loss, potentially leading to improve the accuracy of our 1480 equation.

1481 However, the empirical results, as depicted in Fig. 29, demonstrate that the fitting improvement with 1482 the inclusion of ζ is quite marginal when compared to the version without this parameter. This 1483 slight enhancement does not justify the additional complexity introduced by managing negative values of S_2 . Furthermore, the empirical observation that ζ tends to converge close to 1 (e.g. 1.125) 1484 in Fig. 29c) reinforces the notion that the original formulation of the function, without the power 1485 term ζ , is adequately robust. This finding suggests that the direct influence of the learning rate 1486 annealing area, as initially modeled, sufficiently captures the essential dynamics without the need 1487 for this additional complexity. Another additional complexity lies in that S_2^{ζ} becomes incalculable 1488 when $S_2 < 0$ in LR re-warmup. 1489

Studies of scaling laws are mostly empirically driven. Over-parameterizing the scaling law equation 1490 essentially leads to more accurate fitting results. However, it will also complicate the final format 1491 and hinders us to focus on major factors for the training dynamics. We choose our main format not 1492 due to absolute prediction accuracy but to pursue the simplification (i.e., fewest extra parameters) 1493 to model the essential training dynamics of LLMs. Notably, in our main format in Eq. 1, we only 1494 introduce one extra parameter compared to Chinchilla scaling law, i.e., the coefficient C of the S_2 1495 term. As suggested in Sec. 3, our main scaling law format still has a strong and robust capacity 1496 across many practical scenarios. We believe and expect that there should be more powerful specific 1497 format (maybe with more parameters) after this work. 1498

1490 1499 1500

1511

H.4 OPTIMIZING LEARNING RATE SCHEDULE

1501A natural next step of this work would be optimizing LR schedule based on our proposed scaling law.1502From a practical engineering aspect, it is feasible and efficient to select better LRS from many candi-1503dates based on the prediction of the scaling law. Specifically, WSD (rather than cosine) LRS should1504be used to confirm the larger values for both S_1 and S_2 , as stated in Sec. 4.2 and Appendix G.2;1505WSD LRS annealing ratio can be determined by the method stated in Sec. 4.3; Annealing function1506can be selected by the method stated in Appendix G.3. It is quite easy to get the (nearly) optimal LR1507schedule based on the composition of the methods above.

From another perspective, one might adopt our scaling law to directly solve optimal LR schedule mathematically. It could be found that our Eq. 1 leads to a collapsed LRS: some zero learning rates at last. This problem is related to the issue ($\eta \approx 0$ case) that we discussed above in Appendix H.3.

⁷We still recommend future developers to try this format if possible.

Mathematical optimization strongly depends on an absolute accuracy, while our scaling law in such scenario does not achieve perfectly 100% accuracy (shown in Fig. 3). In comparison, our mentioned variant forms (e.g., Eq. 10) with extra parameters should be more preferably used to solve the op-timal LRS, because they cover more corner cases and have higher accuracy. We believe that in the future, there will be stronger and more parameterized specific forms, which are more suitable for directly mathematically solving the optimal LR schedule. At this stage, we believe that the approach based on the practical engineering mentioned above is sufficient to obtain a (nearly) optimal LR schedule.