

# How to Set the Learning Rate for Large-Scale Pre-training?

Anonymous ACL submission

## Abstract

Optimal configuration of the learning rate (LR) is a fundamental yet formidable challenge in large-scale pre-training. Given the stringent trade-off between training costs and model performance, the pivotal question is whether the optimal LR can be accurately extrapolated from low-cost experiments. In this paper, we formalize this investigation into two distinct research paradigms: Fitting and Transfer. Within the Fitting Paradigm, we innovatively introduce a Scaling Law for search factor, effectively reducing the search complexity from  $\mathcal{O}(n^3)$  to  $\mathcal{O}(n \cdot C_D \cdot C_\eta)$  via predictive modeling. Within the Transfer Paradigm, we extend the principles of  $\mu$ Transfer to the Mixture of Experts (MoE) architecture, broadening its applicability to encompass model depth, weight decay, and token horizons.

By pushing the boundaries of existing hyperparameter research in terms of scale, we conduct a comprehensive comparison between these two paradigms. Our empirical results challenge the scalability of the widely adopted  $\mu$ Transfer in large-scale pre-training scenarios. Furthermore, we provide a rigorous analysis through the dual lenses of training stability and feature learning to elucidate the underlying reasons why module-wise parameter tuning underperforms in large-scale settings. This work offers systematic practical guidelines and a fresh theoretical perspective for optimizing industrial-level pre-training.

## 1 Introduction

The rapid evolution of Large Language Models (LLMs) (OpenAI et al., 2024c,a,b, 2025; Team et al., 2025a, 2024; Comanici et al., 2025; DeepSeek-AI et al., 2024a,b, 2025b,a,c) is continuously pushing the cognitive boundaries of artificial intelligence, driven fundamentally by the Scaling Laws (Kaplan et al., 2020) arising from large-scale pre-training. However, executing such large-scale

pre-training remains formidable. A fundamental challenge is selecting an appropriate/optimal learning rate (LR). On one hand, large-scale pre-training involves massive computational loads and prolonged training cycles, requiring a precise LR to ensure both stability and convergence efficiency. On any other hand, the vast consumption of computational resources makes the cost of trial-and-error unacceptable. Consequently, **the crux of learning rate for large-scale pre-training lies in accurately characterizing the relationship between the optimal LR in “cheaper-to-train” small-scale experiments and that of the target scale.**

This paper establishes two fundamental research paradigms for setting the learning rate in large-scale pre-training: **Fitting** and **Transfer**. The Fitting Paradigm involves directly modeling the relationship between the optimal learning rate, model size, and training data under standard initialization conditions, thereby extrapolating the learning rate for the target training scale (DeepSeek-AI et al., 2024a; Li et al., 2025). To overcome the bottlenecks of combinatorial explosion and prohibitive training costs inherent in prior research within the fitting paradigm, this work innovatively introduces a scaling Law for search factor. By leveraging performance prediction, we effectively reduce the search complexity from  $\mathcal{O}(n^3)$  to  $\mathcal{O}(n \cdot C_D \cdot C_\eta)$ .

The Transfer Paradigm, on the other hand, conducts hyperparameter optimization (including learning rate) on selected proxy models and subsequently transfers these hyperparameters to the target model according to established rules. In this study, we adopt the fundamental principles of  $\mu$ Transfer (Yang et al., 2022). However, to better align with contemporary large-scale pre-training scenarios, we implement a critical extension by selecting the Mixture of Experts (MoE) as our research architecture. Building upon existing literature, we expand the transfer dimensions to encom-

084 pass model widths and depths, while simultane- 134  
085 ously incorporating the influences of weight decay 135  
086 and token horizon on  $\mu$ Transfer. These enhance- 136  
087 ments substantially push the boundary of applica- 137  
088 bility for  $\mu$ Transfer. 138

089 To facilitate a comprehensive comparison 139  
090 between the two paradigms, this study extends 140  
091 the target prediction scale for learning rates by 141  
092 more than tenfold (x10). Our target configura- 142  
093 tion is set as a MoE model with 12B total param- 143  
094 eters with 1.3B activated for each token, trained 144  
095 on 500B tokens—a scale that significantly sur- 145  
096 passes existing hyperparameter research. The 146  
097 primary contributions of this paper can be sum- 147  
098 marized in the following three aspects: 148

099 **Paradigm and Theoretical Innovation:** We sys- 149  
100 tematically formalize the two research paradigms 150  
101 and innovatively integrate Scaling Laws for 151  
102 performance prediction. This approach effectively 152  
103 reduces modeling costs while substantially 153  
104 enhancing both prediction efficiency and the range 154  
105 of parameter coverage. 155

106 **Ultra-Large-Scale Empirical Comparison:** 156  
107 Breaking through the scale limitations of prior 157  
108 hyperparameter studies, this work provides the first 158  
109 comprehensive comparison of the two paradigms 159  
110 within a real-world, large-scale pre-training envi- 160  
111 ronment, offering systematic practical guidelines 161  
112 for large-model engineering. 162

113 **Multidimensional Mechanistic Insights:** We 163  
114 provide an in-depth analysis of the dynamical 164  
115 characteristics of both paradigms during pre-training, 165  
116 focusing on two core dimensions: Training 166  
117 Stability and Feature Learning. This offers a 167  
118 novel perspective for research into large-scale 168  
119 pre-training. 169  
120

## 121 2 Related Works 170

### 122 2.1 Learning Rate Schedule 171

123 Prior to the advent of large-scale language model 172  
124 (LLM) pre-training, the cosine annealing schedule 173  
125 (Loshchilov and Hutter, 2017) served as the pre- 174  
126 dominant standard. However, the cosine schedule 175  
127 mandates a predetermined number of total training 176  
128 steps, rendering it insufficiently flexible amidst the 177  
129 backdrop of continuously expanding pre-training 178  
130 scales. Consequently, the Warmup-Stable-Decay 179  
131 (WSD) schedule (Hu et al., 2024) has emerged. 180  
132 This schedule is characterized by a stable phase 181  
133 where the learning rate remains constant follow-

ing the warmup period, eventually decaying to a 134  
specific terminal value. Since the decay phase can 135  
be initiated at any point during the stable phase 136  
to conclude training, WSD is regarded as highly 137  
adaptable to the dynamic requirements of large- 138  
scale pre-training. Reflecting this advantage, the 139  
WSD scheduler has recently been adopted by main- 140  
stream large-scale pre-training projects(DeepSeek- 141  
AI et al., 2025b; Team et al., 2025b; Bai et al., 142  
2025). 143

Propelled by scaling laws, the magnitude of pre- 144  
training continues to escalate. The stable phase 145  
frequently spans weeks or even months (DeepSeek- 146  
AI et al., 2025b; Bai et al., 2025; Yang et al., 2025), 147  
making the precise configuration of the learning 148  
rate critically important. However, existing re- 149  
search on learning rate configuration has predom- 150  
inantly focused on the cosine annealing schedule. 151  
Under the cosine regime, Kaplan et al. (2020) elu- 152  
cidated the relationship between the learning rate and 153  
model parameters, while Bjorck et al. (2025) and 154  
Li et al. (2025) empirically derived power-law for- 155  
mulations correlating the learning rate with model 156  
size  $N$  and training data size  $D$ . Diverging from 157  
existing literature, our work investigates the rela- 158  
tionship between the optimal learning rate, model 159  
size, and training data size specifically within the 160  
stable phase of a constant learning rate schedule. 161

### 162 2.2 Maximal Update Parametrization 162

Maximal Update Parametrization 163  
( $\mu$ Parametrization or  $\mu$ P, Yang et al. (2022)) 164  
is a widely investigated framework for hyperpa- 165  
rameter configuration. The fundamental premise 166  
of  $\mu$ P is to guarantee training stability and 167  
ensure that weights across different modules are 168  
adequately trained(i.e. maximal feature learning) 169  
even as model width approaches infinity. 170

By virtue of maintaining these properties in the 171  
infinite-width limit,  $\mu$ P possesses inherent capabil- 172  
ities for hyperparameter transfer. This gives rise to 173  
a derivative method known as  $\mu$ Transfer, wherein 174  
the optimal learning rate for a target model can be 175  
directly calculated based on the optimum identi- 176  
fied via search on a smaller proxy model. While 177  
the initial formulation of  $\mu$ Transfer was limited 178  
to extrapolating model width, subsequent stud- 179  
ies by Yang et al. (2023) and Dey et al. (2025) 180  
have investigated extensions for scaling model 181  
depth. Beyond its extensive application in dense 182  
architectures(Lingle, 2025),  $\mu$ Transfer has also 183  
been experimentally applied to Mixture-of-Experts 184

(MoE) structures(Małaśnicki et al., 2025). Furthermore, recent research indicates that the efficacy of  $\mu$ Transfer is primarily manifested during the early stages of training; to extend the effective transfer horizon, adjustments to weight decay are required(Wang and Aitchison, 2025; Małaśnicki et al., 2025; Fan et al., 2025).

Building upon existing research of  $\mu$ Transfer and integrating current methodologies for pre-training hyperparameter configuration, our work conducts a granular investigation into the impact of  $\mu$ Transfer on the performance of large-scale pre-training.

### 3 Approach

This section delineates the specific methodologies for the configuration of the learning rate under two distinct paradigms. Section 3.1 introduces the Fitting Paradigm, which leverages scaling laws to enhance the efficiency and scope of the fitting process. Section 3.2 focuses on the representative transfer paradigm  $\mu$ Transfer method and elucidating its practical implementation within large-scale pre-training contexts. Crucially, our study focuses on the stable training phase governed by the Warmup-Stable-Decay (WSD) learning rate schedule.

#### 3.1 Scaling Laws for Learning Rate

For a given model size  $N$  and training data size  $D$ , the optimal learning rate  $\eta$  is formulated as:

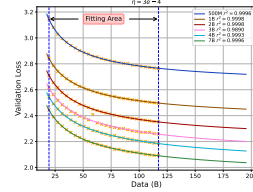
$$\eta_{ND}^* = \underset{\eta}{\operatorname{argmin}} L(\eta | N, D, \Theta), \quad (1)$$

where  $L$  is validation loss and  $\Theta$  contains other hyperparameters involved in the pre-train process.

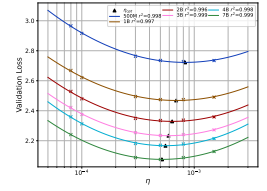
Characterizing the relationship between the optimal  $\eta$ , model size, and data size necessitates a grid search across the  $N$ ,  $D$ , and  $\eta$  dimensions, resulting in a computational complexity of approximately  $\mathcal{O}(n^3)$ . Fortunately, inspired by prior work (Bjorck et al., 2025), we observe that for fixed  $N$  and  $D$ , the relationship between the validation loss  $L$  and the learning rate  $\eta$  approximates an invex profile. Consequently, we employ a quadratic polynomial to fit this relationship:

$$L(\eta | N, D, \Theta) = L_{min} + C \cdot (\log(\eta) - \eta_{min})^2, \quad (2)$$

where  $L_{min}$ ,  $C$ , and  $\eta_{min}$  are the fitting coefficients.



(a) Fitting results for Equation 4. Data points to the left of the dashed line represent the empirical values used for fitting; The curves to the right depicts predictions. See Appendix A.2 for the discussion on the accuracy of Equation 4.



(b) Results of fitting the validation loss against the learning rate (LR) using a **quadratic polynomial**. Different colored curves correspond to models of varying sizes, while the triangle indicate the optimal LR.

Figure 1: Results of Equation 4 and 2. These approaches allow for a substantial reduction in the time and storage cost of the search process.

Consequently, for a given  $N$  and  $D$ , the optimal learning rate  $\eta^*$  can be directly derived via fitting on a limited set of learning rates:

$$\log(\eta^*) = \eta_{min} = \underset{\eta}{\operatorname{argmin}} \{L(\eta | N, D, \Theta)\}, \quad (3)$$

Figure 1(b) shows the fitted curves of Equation 2. The search complexity is reduced from  $\mathcal{O}(n^3)$  to  $\mathcal{O}(n^2 * C_\eta)$

Furthermore, inspired by contemporary research on scaling laws(Hoffmann et al., 2022; Tissue et al., 2024), we observe that under the WSD schedule, the validation loss exhibits a power-law relationship with the training data size  $D$  for a fixed model configuration:

$$L(D) = L_0 + A \cdot D^{-\gamma}, \quad (4)$$

Where  $L_0$ ,  $A$ ,  $\gamma$  are parameters to fit. This implies that the search space along the dimension of data size  $D$  can be significantly compressed, enabling the extrapolation of results to larger data regimes via a limited number of search points. The specific fitting procedure is illustrated in Figure 1(a). This methodology effectively improves the trade-off between search cost and fitting precision, reducing the computational complexity of the search from  $\mathcal{O}(n^3)$  to  $\mathcal{O}(n \cdot C_D \cdot C_\eta)$ .

Based on Equation 2 and 3, by conducting a search within the  $N$  dimension, we can efficiently derive a comprehensive set of optimal LR  $\{\eta_{ND}^*\}_{N,D}$ , corresponding to varying model sizes  $N$  and data sizes  $D$ . This facilitates the fitting of the functional relationship between the optimal LR and the variables  $N$  and  $D$ :

$$Lr(N, D) = \operatorname{argmin}_{Lr \in \mathcal{F}} L(Lr(N, D), \eta_{ND}^* \mid \Theta), \quad (5)$$

where  $\mathcal{F}$  represents the candidate function space, and  $L$  denotes the metric function, which is Root Mean Squared Error (RMSE) in our work.

The final fitted relationship governing the optimal learning rate with respect to model size  $N$  and data size  $D$  is given by:

$$Lr(N, D) = 38.4588 \cdot N^{-0.2219} \cdot D^{-0.3509}. \quad (6)$$

We observe a good fit with  $R^2 \approx 0.9622$  (See Appendix A.3.1 for details of fitting process). The overall fitting results are shown in Figure 2.

Extending this approach, we further conduct a fine-grained investigation into the learning rate configurations for distinct model modules in Section 6.1.

### 3.2 Scaling $\mu$ Transfer for Pre-training

As the Mixture-of-Experts (MoE) architecture increasingly serves as the foundational backbone for large-scale pre-training (DeepSeek-AI et al., 2024b, 2025b,a,c; Yang et al., 2024; Qwen et al., 2025; Yang et al., 2025; Bai et al., 2025), we adopt the MoE architecture as our proxy model for  $\mu$ P. Regarding the target model, we adhere to the settings proposed by Małański et al. (2025) for initialization along the width dimension. For the depth dimension, we draw upon the methodologies of Depth-up (Yang et al., 2023) and Complete- $\mu$ P (Dey et al., 2025; Mlodozieniec et al., 2025). The central mechanism involves applying a depth-dependent scaling factor to the residual branch:

$$H^{i+1} = H^i + m_L^{-\alpha} \mathcal{F}(H^i), \quad i \in \{1, \dots, L\}, \quad (7)$$

where  $H^i$  denotes the output of the  $i$ -th layer, and  $\mathcal{F}$  represents either the Attention or Feed-Forward Network (FFN) layer. Following the recommendations of Complete- $\mu$ P, we set  $\alpha = 1$  to enhance the transferability of  $\mu$ Transfer.

Wang and Aitchison (2025) and Fan et al. (2025) have identified weight decay  $\lambda$  as a critical determinant of  $\mu$ Transfer efficacy. Consequently, we incorporate the influence of weight decay into the training process of the target model, maintaining the proportionality  $\delta\lambda \propto \delta lr$ . For given model

size  $N$  and data volume  $D$ , we observe that the approximate invex relationship between validation loss  $L$  and learning rate persists within  $\mu$ P proxy models. This observation allows for a reduction in the search space along the learning rate dimension, thereby improving the efficiency of  $\mu$ Transfer. Regarding transfer along the token horizon dimension, we adopt the configuration from Mlodozieniec et al. (2025). The detailed initialization and transfer rules for the target model parameters are summarized in Table 2 and Table 10.

## 4 Experiments

### 4.1 Datasets

The pre-training corpus utilized in our work is derived from InternLM2.5 (Cai et al., 2024), including general text, source code, and long-context sequences. Specifically, the textual component spans web pages, academic papers, patents, and books. The code component is primarily sourced from GitHub, programming communities, and other public repositories, covering a diverse array of programming languages including C/C++, Java, and Python. All data underwent rigorous deduplication and safety filtering protocols.

To ensure distributional consistency, the validation set employed in our experiments was constructed via random sampling from the above corpus, while strictly maintaining disjointness from the training samples to prevent data leakage.

### 4.2 Experimental Settings

We adopt the Qwen3-MoE (Yang et al., 2025) architecture for our experimental models. For all model training, we utilize the AdamW optimizer (Loshchilov and Hutter, 2019) with  $\beta_1 = 0.9$ ,  $\beta_2 = 0.95$ , and  $\epsilon = 10^{-8}$ . The learning rate schedule consists of a linear warmup for 1,000 steps, followed by a constant learning rate strategy. The sequence length is fixed at 4,096, and the global batch size is set to 4M tokens.

For experiments in 3.1, we employ models of four distinct sizes (550M, 1B, 2B, and 3B) all adhering to the structural configuration of the Qwen3-30B-A3B model. Notably, the aspect ratio between model width and depth remains constant across these scales. We subsequently validate our experimental findings on target models with 4B, 12B total parameters. With the exception of normalization parameters, all model weights are initialized from a normal distribution with a standard devia-

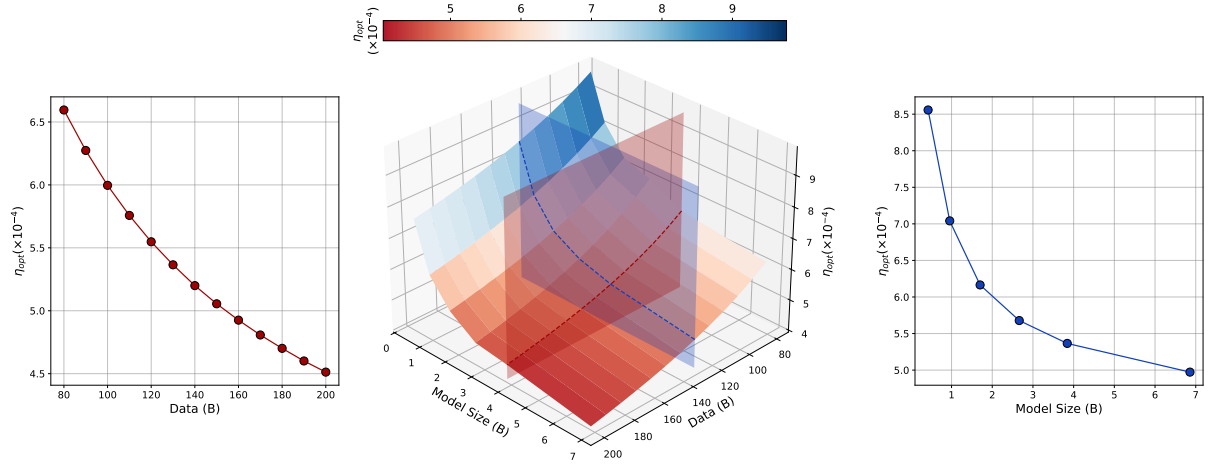


Figure 2: **Middle**: Visualization of the optimal learning rate relative to model size  $N$  and data size  $D$ . **Left**: The relationship between the optimal learning rate and data size  $D$  with model size fixed at  $N = 4B$ . **Right**: The relationship between the optimal learning rate and model size  $N$  with data size fixed at  $D = 140B$ .

tion of 0.02. The search space for the learning rate is defined as  $\eta \in \{8e - 5, 1e - 4, 3e - 4, 5e - 4, 8e - 4, 1.5e - 3, 2e - 3\}$ . Each model is trained on approximately 120B tokens (30,000 steps), and we extrapolate the results to 500B tokens using Equation 4. Weight decay is set to 0.1.

For the  $\mu$ Transfer experiments, we conduct learning rate and initialization searches on a proxy model with 2B total parameters. The search space for the learning rate is defined as  $\eta \in \{8e - 5, 1e - 4, 3e - 4, 5e - 4, 8e - 4, 1.5e - 3, 2e - 3\}$ ; for initialization, we explore the range  $\sigma \in \{0.0005, 0.001, 0.002, 0.005, 0.01, 0.015, 0.02\}$ . The actual training data size for these experiments is approximately 200B tokens (50,000 steps), and we extrapolate the hyperparameters to a 500B token regime with Equation 4 and settings from Mlodozienec et al. (2025).

### 4.3 Evaluation

To assess the downstream performance of the models developed during our validation experiments, we evaluate our models on MMLU (Hendrycks et al., 2021) and CMMLU (Li et al., 2024) benchmarks. MMLU serves as our primary English evaluation set, comprising four-choice multiple-choice questions across 57 distinct subjects, including anatomy, physics, genetics etc. Conversely, we employ CMMLU to evaluate Chinese language proficiency which covers 67 domains ranging from natural sciences and humanities to general knowledge.

For the implementation of these evaluations, we leverage the OpenCompass framework (Contribu-

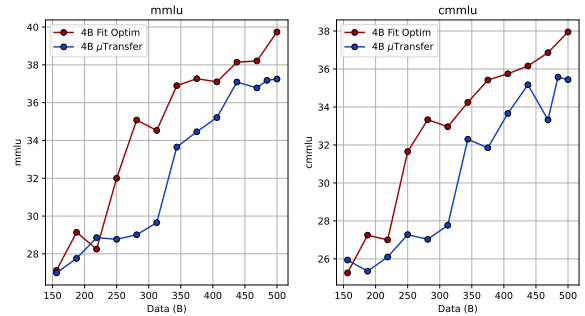


Figure 3: Downstream task performances of 4B model with global optimal LR and  $\mu$ P respectively.

tors, 2023b), a comprehensive Python toolkit designed to facilitate the batch evaluation of diverse foundation models across heterogeneous datasets. Furthermore, to expedite the evaluation pipeline, we utilize the LMDeploy framework (Contributors, 2023a) with Turbomind (Zhang et al., 2025) backend for efficient model loading and inference acceleration.

## 5 Results

First, we extrapolate the proxy model solely by increasing its width, scaling it to 4B total parameters with 530M active parameters, and conducting from-scratch pre-training on 500B tokens. To rigorously assess the pre-training quality under both paradigms, we evaluate not only the final model performance but also the downstream task results throughout the training process. The performance trends are illustrated in Figure 3. As shown, the pre-training quality achieved by the Fitting Paradigm

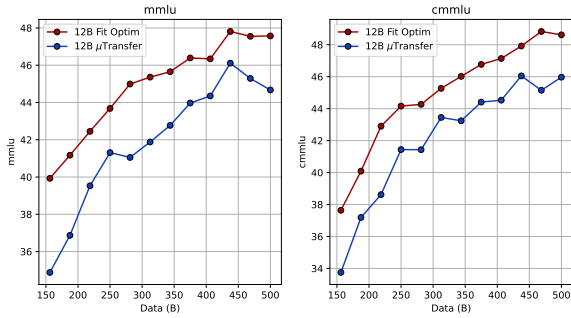


Figure 4: Downstream task performances of 12B model with global optimal LR and  $\mu P$  respectively.

is significantly higher than that of  $\mu Transfer$ .

Furthermore, we extend the predictive training scale by more than an order of magnitude, scaling the model to 12B total parameters (1.3B active) and pre-training it from scratch on 500B tokens. We similarly evaluate the intermediate progress, with the overall performance trajectories presented in Figure 4. As demonstrated in Figure 4, the model trained via the Fitting Paradigm consistently and significantly outperforms the one using  $\mu Transfer$ .

## 6 Analyze

### 6.1 Module-Level Optimal Learning Rates

A fundamental motivation behind  $\mu P$  is the hypothesis that under Standard Parametrization (SP) and a uniform global learning rate, specific modules may suffer from insufficient training, thereby failing to satisfy the regime of maximal feature learning. To investigate this, building upon the global optimal learning rate derived from our fitting paradigm in Section 3.1, we employ a greedy search strategy to conduct a fine-grained learning rate search across four distinct parameter modules: Embeddings, LM Head, Router, and Hidden parameters. We observe that fine-grained tuning of individual modules yields no significant performance improvement compared to the global optimal learning rate configuration.

The optimal learning rates identified for specific modules align closely with the global optimum, and the minimum loss achieved through module-specific search exhibits negligible deviation from the loss achieved with global optimal learning rate from Equation 6 (as illustrated in Figure 5). Consequently, assigning distinct optimal learning rates to specific modules does not appear to materially enhance model performance.

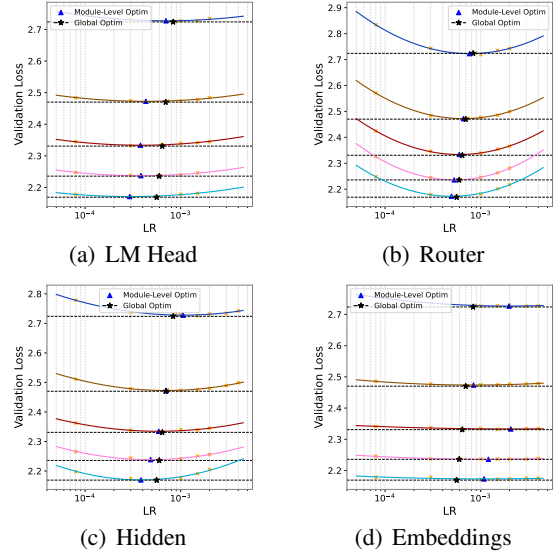


Figure 5: The relationship between loss and learning rate for (a)LM Head, (b)Router, (c)Hidden and (d)Embedding parameters during the module-level learning rate search. Each curve corresponds to a model of a specific size. The dashed lines indicate the loss achieved by the corresponding model size under the global optimal learning rate setting. Triangle markers denote the optimal learning rate for the current module, while star markers represent the global optimal learning rate.

To further validate the effect of module-level optimal LR, we trained the target 4B model for 120B tokens using both the derived module-specific optimal LR and the calculated global optimal LR. As depicted in Figure 6, the comparison of the validation losses reveals that the loss curves for both settings are virtually indistinguishable. See Table 3 and Appendix A.3.2 for detailed settings.

### 6.2 A Closer Look at Feature Learning

In the previous subsection, we observed that fine-grained learning rate tuning across distinct model modules yielded no substantial performance gains, indicating that a global learning rate configuration does not induce training imbalances among components. In this subsection, we further investigate the feature learning dynamics of these modules by analyzing the optimization trajectory, specifically monitoring the evolution of parameter update magnitudes throughout the training process.

As illustrated in the Figure 7, training with the AdamW optimizer results in parameter update magnitudes that remain stable over extended periods and exhibit relative uniformity across layers. The update magnitudes consistently approximate 0.2, a finding consistent with recent theoretical studies

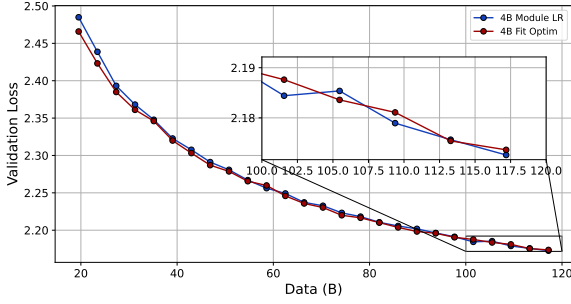


Figure 6: Performance comparison between the global optimal LR (red line) and module-wise optimal LR (blue line) on a 4B model trained for 120B tokens. The two loss curves are virtually indistinguishable during the mid-to-late stages of training ( $\Delta L \leq 0.01$ ), indicating that module-specific learning rate optimization does not yield significant performance improvements.

(Liu et al., 2025; Kosson et al., 2024). This evidence further corroborates that distinct modules maintain comparable feature learning capabilities at any given stage of training, thereby negating the necessity for module-specific learning rates to balance feature learning efficiency.

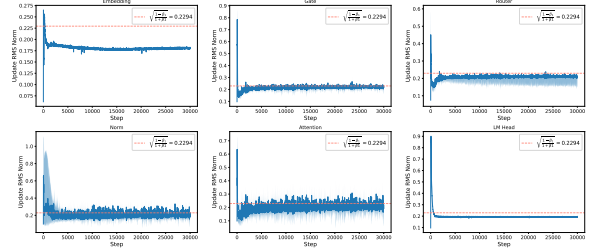
### 6.3 Does Standard Parametrization Scale?

Training stability is widely regarded as another distinct advantage of  $\mu P$ . Yang et al. (2022) argues that under standard parametrization, the internal training states of certain modules tend to "blow up" as model scale increases, thus the adjustment of learning rates on different modules is necessary.

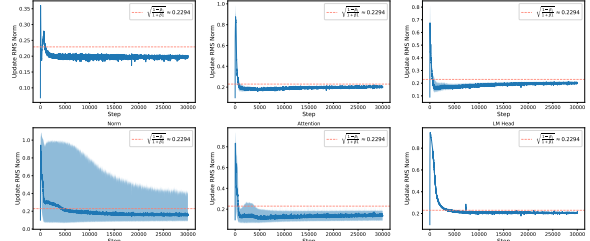
We replicated the methodology of  $\mu$ Transfer to analyze the model derived from our experiments. Contrary to expectations, under standard parametrization, the internal states of our model did not exhibit blow up (Figure 8(a)); rather, they displayed trends remarkably similar to those of models initialized via  $\mu P$ . To investigate further, we conducted an ablation where the QK-Norm modules were removed during the computation of attention logits (Figure 9). Under this condition, we successfully reproduced the instability trends described in  $\mu$ Transfer. Consequently, we posit that recent advancements in model architecture—such as the incorporation of QK-Norm—have rendered layer-wise training more balanced and significantly enhanced robustness to hyperparameter variations.

### 6.4 Impact of Data Size on Training Stability

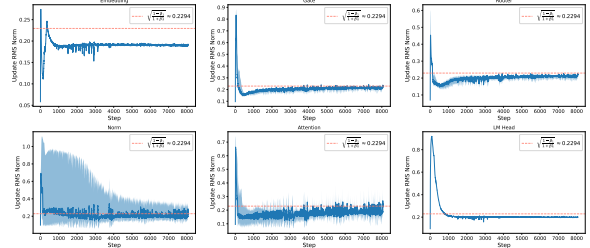
Existing research on training stability, most notably the work on  $\mu P$ , has predominantly focused on model scale while neglecting the influence of



(a) 4B Model

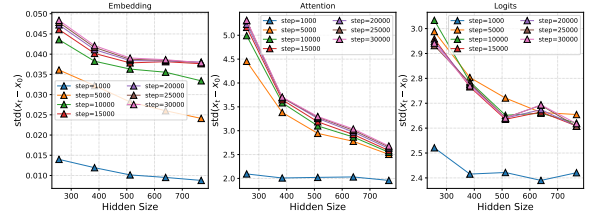


(b) 500M Model without QK-Norm

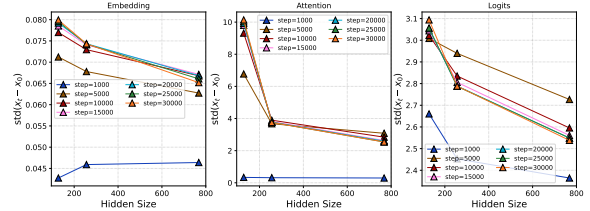


(c) 4B  $\mu P$  Model

Figure 7: Update RMS norm during the training process of: (a) 4B model (b) 500M model with QK-Norm removed (c) 4B model with  $\mu P$ . Update RMS norm maintains approximately 0.2 in all the models we observed.



(a) 4B Model



(b) 4B  $\mu P$  Model

Figure 8: Variation of word embeddings, attention logits, and logits compared to initial states at certain training steps as width increases. With reference to Yang et al. (2022), we plot the standard deviation of the coordinates of  $x_t - x_0$ ,  $x \in \{word\ embeddings, attention\ logits, logits\}$ . In our experiments, logits and attention logits of models with standard parametrization do not exhibit the "blow-up" tendency.

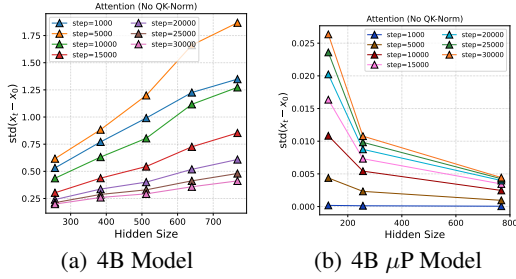


Figure 9: Variation of attention logits at certain training steps as width increases. We ignored QK-Norm parameters when compute attention logits. Attention logits started to blow up with width in SP model.

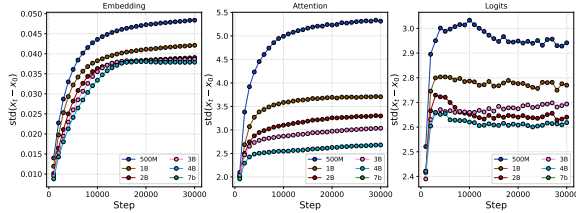


Figure 10: Variation of word embeddings, attention logits, and logits compared to initial states as training proceeded.

training data size. Employing the analytical framework established in Section 6.3, we investigate the evolution of the model’s internal states as the amount of training data increases under standard parametrization.

As illustrated in the Figure 10, while the model’s internal states eventually converge to a relatively stable regime as the amount of training data increases, the internal variations are significantly more drastic with respect to data scaling than to model scaling. This phenomenon is particularly evident in the attention logits. This observation offers a potential explanation for the scaling coefficients in Equation 6, where the exponent for model parameter count  $N$  ( $-0.22$ ) is algebraically greater than that for data volume  $D$  ( $-0.35$ ). As the size of training data expands, the magnitude of parameter updates across different modules exhibits a more pronounced increase; consequently, the optimal learning rate requires more substantial adjustment to maintain training stability.

## 6.5 Decay Training

A key characteristic of WSD schedule is the utilization of higher-quality training data during the decay phase after the constant-learning-rate stable phase to maximize the model’s feature learning.

Building upon the experiments described in Section 5 we extended the training of both model vari-

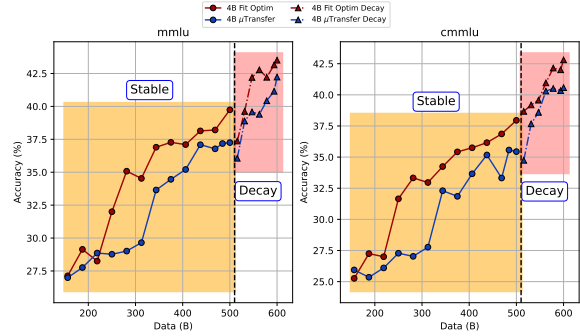


Figure 11: Downstream task performances of 4B model after the end of decay phase. The left area depicts the stable phase of training, while the right area corresponds to the decay phase.

ants after the end of the stable phase using a distinct corpus of high-quality data. We annealed the learning rate to 10% of its value during the stable phase and continued training for an additional 100B tokens. We then evaluated the downstream task performance of the models that had completed the full WSD training.

As shown in Figure 11, the model trained with the global optimal learning rate outperformed the model derived from  $\mu$  Transfer on both the MMLU and CMMLU benchmarks, achieving accuracy improvements of 1.28% (42.23%  $\rightarrow$  43.51%) and 2.23% (40.58%  $\rightarrow$  42.81%), respectively. These results demonstrate that the global optimal learning rate yields superior performance in realistic pre-training scenarios.

## 7 Conclusion

This paper systematically establishes two fundamental research paradigms—Fitting and Transfer—to address the critical challenge of learning rate configuration in large-scale pre-training. At the methodological level, we introduce scaling Laws to reduce the complexity of the Fitting Paradigm, and provide a comprehensive extension of  $\mu$ Transfer across model architectures, depths, weight decay, and token horizons. Through extensive experimentation, we challenge the scalability of the widely adopted  $\mu$ Transfer in large-scale pre-training scenarios and provide an in-depth analysis of the underlying mechanisms that limit the performance of module-wise parameter tuning at scale. This research offers both systematic practical guidance and a novel theoretical perspective for optimizing industrial-level pre-training.

## 561 Limitations

562 To inform and inspire future research, we summa-  
563 rize the limitations of our work as follows:

564 Learning Rate Schedules: This study focuses on  
565 large-scale pre-training, where the Warm-Stable-  
566 Decay (WSD) scheduler is currently the industry  
567 standard. Consequently, our analysis is centered  
568 on this specific schedule and does not explore the  
569 dynamics of other learning rate schedulers.

570 Model Architectures: Given that the Mixture  
571 of Experts (MoE) architecture has become the  
572 foundational backbone for modern large-scale  
573 pre-training, it served as the primary subject of our  
574 investigation. The generalizability of our findings  
575 to Dense architectures remains to be verified in  
576 future work.

577 Extrapolation Limits: Due to computational re-  
578 source constraints, this study did not investigate the  
579 ultimate extrapolation boundaries (i.e., the maxi-  
580 mum scale at which these predictions remain accu-  
581 rate) for both the Fitting and Transfer paradigms.

## 585 Use of AI Assistants

586 We primarily use AI assistants to improve and en-  
587 rich our writing.

## 588 References

589 Lei Bai, Zhongrui Cai, Yuhang Cao, Maosong Cao,  
590 Weihan Cao, Chiyu Chen, Haojiong Chen, Kai  
591 Chen, Pengcheng Chen, Ying Chen, Yongkang Chen,  
592 Yu Cheng, Pei Chu, Tao Chu, Erfei Cui, Ganqu Cui,  
593 Long Cui, Ziyun Cui, Nianchen Deng, and 158 oth-  
594 ers. 2025. *Intern-s1: A scientific multimodal founda-  
595 tion model*. *Preprint*, arXiv:2508.15763.

596 Johan Bjorck, Alon Benhaim, Vishrav Chaudhary, Furu  
597 Wei, and Xia Song. 2025. *Scaling optimal lr across  
598 token horizons*. *Preprint*, arXiv:2409.19913.

599 Zheng Cai, Maosong Cao, Haojiong Chen, Kai Chen,  
600 Keyu Chen, Xin Chen, Xun Chen, Zehui Chen, Zhi  
601 Chen, Pei Chu, Xiaoyi Dong, Haodong Duan, Qi Fan,  
602 Zhaoye Fei, Yang Gao, Jiaye Ge, Chenya Gu, Yuzhe  
603 Gu, Tao Gui, and 81 others. 2024. *Internlm2 techni-  
604 cal report*. *Preprint*, arXiv:2403.17297.

605 Gheorghe Comanici, Eric Bieber, Mike Schaekermann,  
606 Ice Pasupat, Noveen Sachdeva, Inderjit Dhillon, Mar-  
607 cel Blistein, Ori Ram, Dan Zhang, Evan Rosen, Luke  
608 Marris, Sam Petulla, Colin Gaffney, Asaf Aharoni,  
609 Nathan Lintz, Tiago Cardal Pais, Henrik Jacobs-  
610 son, Idan Szpektor, Nan-Jiang Jiang, and 3416 oth-  
611 ers. 2025. *Gemini 2.5: Pushing the frontier with*

advanced reasoning, multimodality, long context,  
and next generation agentic capabilities. *Preprint*,  
arXiv:2507.06261. 612  
613  
614

LMDeploy Contributors. 2023a. Lmdeploy: A toolkit  
for compressing, deploying, and serving llm. <https://github.com/InternLM/lmdeploy>. 615  
616  
617

OpenCompass Contributors. 2023b. Opencompass:  
A universal evaluation platform for foundation  
models. [https://github.com/open-compass/  
opencompass](https://github.com/open-compass/opencompass). 618  
619  
620  
621

DeepSeek-AI, :, Xiao Bi, Deli Chen, Guanting Chen,  
Shanhuang Chen, Damai Dai, Chengqi Deng,  
Honghui Ding, Kai Dong, Qishui Du, Zhe Fu,  
Huazuo Gao, Kaige Gao, Wenjun Gao, Ruiqi  
Ge, Kang Guan, Daya Guo, Jianzhong Guo, and  
69 others. 2024a. *Deepseek llm: Scaling open-  
source language models with longtermism*. *Preprint*,  
arXiv:2401.02954. 622  
623  
624  
625  
626  
627  
628  
629

DeepSeek-AI, Daya Guo, Dejian Yang, Haowei Zhang,  
Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu,  
Shirong Ma, Peiyi Wang, Xiao Bi, Xiaokang Zhang,  
Xingkai Yu, Yu Wu, Z. F. Wu, Zhibin Gou, Zhi-  
hong Shao, Zhuoshu Li, Ziyi Gao, and 181 others.  
2025a. *Deepseek-r1: Incentivizing reasoning capa-  
bility in llms via reinforcement learning*. *Preprint*,  
arXiv:2501.12948. 630  
631  
632  
633  
634  
635  
636  
637

DeepSeek-AI, Aixin Liu, Bei Feng, Bin Wang, Bingx-  
uan Wang, Bo Liu, Chenggang Zhao, Chengqi Deng,  
Chong Ruan, Damai Dai, Daya Guo, Dejian Yang,  
Deli Chen, Dongjie Ji, Erhang Li, Fangyun Lin, Fuli  
Luo, Guangbo Hao, Guanting Chen, and 138 others.  
2024b. *Deepseek-v2: A strong, economical, and ef-  
ficient mixture-of-experts language model*. *Preprint*,  
arXiv:2405.04434. 638  
639  
640  
641  
642  
643  
644  
645

DeepSeek-AI, Aixin Liu, Bei Feng, Bing Xue, Bingx-  
uan Wang, Bochao Wu, Chengda Lu, Chenggang  
Zhao, Chengqi Deng, Chenyu Zhang, Chong Ruan,  
Damai Dai, Daya Guo, Dejian Yang, Deli Chen,  
Dongjie Ji, Erhang Li, Fangyun Lin, Fucong Dai,  
and 181 others. 2025b. *Deepseek-v3 technical re-  
port*. *Preprint*, arXiv:2412.19437. 646  
647  
648  
649  
650  
651  
652

DeepSeek-AI, Aixin Liu, Aoxue Mei, Bangcai Lin,  
Bing Xue, Bingxuan Wang, Bingzheng Xu, Bochao  
Wu, Bowei Zhang, Chaofan Lin, Chen Dong,  
Chengda Lu, Chenggang Zhao, Chengqi Deng, Chen-  
hao Xu, Chong Ruan, Damai Dai, Daya Guo, De-  
jian Yang, and 245 others. 2025c. *Deepseek-v3.2:  
Pushing the frontier of open large language models*.  
*Preprint*, arXiv:2512.02556. 653  
654  
655  
656  
657  
658  
659  
660

Nolan Dey, Bin Claire Zhang, Lorenzo Noci, Mufan  
Li, Blake Bordelon, Shane Bergsma, Cengiz Pehle-  
van, Boris Hanin, and Joel Hestness. 2025. *Don't  
be lazy: Completex enables compute-efficient deep  
transformers*. *Preprint*, arXiv:2505.01618. 661  
662  
663  
664  
665

Zhiyuan Fan, Yifeng Liu, Qingyue Zhao, Angela Yuan,  
and Quanquan Gu. 2025. *Robust layerwise scal-  
ing rules by proper weight decay tuning*. *Preprint*,  
arXiv:2510.15262. 666  
667  
668  
669



Lazaridou, and 1332 others. 2025a. **Gemini: A family of highly capable multimodal models**. *Preprint*, arXiv:2312.11805.

Gemini Team, Petko Georgiev, Ving Ian Lei, Ryan Burnell, Libin Bai, Anmol Gulati, Garrett Tanzer, Damien Vincent, Zhufeng Pan, Shibo Wang, Soroosh Mariooryad, Yifan Ding, Xinyang Geng, Fred Alcober, Roy Frostig, Mark Omernick, Lexi Walker, Cosmin Paduraru, Christina Sorokin, and 1118 others. 2024. **Gemini 1.5: Unlocking multimodal understanding across millions of tokens of context**. *Preprint*, arXiv:2403.05530.

Kimi Team, Yifan Bai, Yiping Bao, Guanduo Chen, Jiahao Chen, Ningxin Chen, Ruijue Chen, Yanru Chen, Yuankun Chen, Yutian Chen, Zhuofu Chen, Jialei Cui, Hao Ding, Mengnan Dong, Angang Du, Chenzhuang Du, Dikang Du, Yulun Du, Yu Fan, and 150 others. 2025b. **Kimi k2: Open agentic intelligence**. *Preprint*, arXiv:2507.20534.

Howe Tissue, Venus Wang, and Lu Wang. 2024. **Scaling law with learning rate annealing**. *Preprint*, arXiv:2408.11029.

Xi Wang and Laurence Aitchison. 2025. **How to set adam’s weight decay as you scale model and dataset size**. *Preprint*, arXiv:2405.13698.

An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Gao, Chengen Huang, Chenxu Lv, Chujie Zheng, Dayiheng Liu, Fan Zhou, Fei Huang, Feng Hu, Hao Ge, Haoran Wei, Huan Lin, Jialong Tang, and 41 others. 2025. **Qwen3 technical report**. *Preprint*, arXiv:2505.09388.

An Yang, Baosong Yang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Zhou, Chengpeng Li, Chengyuan Li, Dayiheng Liu, Fei Huang, Guanting Dong, Haoran Wei, Huan Lin, Jialong Tang, Jialin Wang, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Ma, and 43 others. 2024. **Qwen2 technical report**. *Preprint*, arXiv:2407.10671.

Greg Yang, Edward J. Hu, Igor Babuschkin, Szymon Sidor, Xiaodong Liu, David Farhi, Nick Ryder, Jakub Pachocki, Weizhu Chen, and Jianfeng Gao. 2022. **Tensor programs v: Tuning large neural networks via zero-shot hyperparameter transfer**. *Preprint*, arXiv:2203.03466.

Greg Yang, Dingli Yu, Chen Zhu, and Soufiane Hayou. 2023. **Tensor programs vi: Feature learning in infinite-depth neural networks**. *Preprint*, arXiv:2310.02244.

Li Zhang, Youhe Jiang, Guoliang He, Xin Chen, Han Lv, Qian Yao, Fangcheng Fu, and Kai Chen. 2025. **Efficient mixed-precision large language model inference with turbomind**. *arXiv preprint arXiv:2508.15601*.

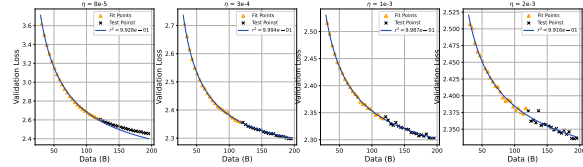


Figure 12: The precision of Equation 4.

## A Appendix

### A.1 Model Architectures and LR Settings

Table 1, 2 and 3 shows the detailed parameters of models’ architectures and learning rate settings.

### A.2 Extrapolation of L(D)

In the experiments presented in our works,  $L(D) = L_0 + A \cdot D^{-\gamma}$  (Equation 4) is repeatedly employed for curve fitting and data augment. To validate the effectiveness of this approach, we continue training the 2B proxy model of  $\mu$ Transfer from 120B to approximately 200B tokens (50,000 steps). Points corresponding to  $\leq 120B$  are used as the fitting set, while the remaining data serve as the test set. The fitted curve is illustrated in Figure 12.

Among the four settings, the extrapolation accuracy is notably poorer under the learning rate of  $8e-5$ , while the predicted values for the other learning rates at 200B tokens align closely with the ground truth. We attribute this discrepancy to the fact that the learning rate of  $8e-5$  is excessively small for the current model, causing the model state to evolve too gradually as data volume increases. By 120B tokens, the model has yet to exhibit a clear trend toward convergence. Consequently, the curve fitted by Equation 4 declines sharply rather than gradually flattening, leading to a misinterpretation of the model’s future trajectory. In contrast, the other tested learning rates are relatively larger and closer to the model’s optimal learning rate, enabling the validation loss curve to enter the convergence phase more rapidly and thus yielding more accurate predictions from the L(D) curve.

For experiments conducted in Section 5, we also conduct the same experiment on 4B Model, fitting with points that corresponding to  $D \leq 120B$ . Figure 13 indicates that despite only approximately one-quarter of the data is used for fitting, the predicted value at  $D = 500B(2.066)$  exhibits a negligible discrepancy from the actual value (2.070). Therefore, we consider the method of data extrapo-

Table 1: Overview of Qwen3-MoE Model Architectures and Hyperparameters in Section 3.1

| Models               | Total Params | Activate Params | Hidden Size | Num Layers | Attn Heads | KV Heads | Interm. Size | Learning Rate                              |
|----------------------|--------------|-----------------|-------------|------------|------------|----------|--------------|--|
| <i>Training Set</i>  |              |                 |             |            |            |          |              |  |
| Qwen3-MoE-0.5B-A0.1B | 550M         | 100M            | 256         | 3          | 32         | 4        | 768          | 8e-5, 1e-4, 3e-4, 5e-4, 8e-4, 1.5e-3, 2e-3 |
| Qwen3-MoE-1B-A0.2B   | 1B           | 190M            | 384         | 9          | 32         | 4        | 768          | 8e-5, 1e-4, 3e-4, 5e-4, 8e-4, 1.5e-3, 2e-3 |
| Qwen3-MoE-2B-A0.3B   | 2B           | 280M            | 512         | 12         | 32         | 4        | 768          | 8e-5, 1e-4, 3e-4, 5e-4, 8e-4, 1.5e-3       |
| Qwen3-MoE-3B-A0.4B   | 3B           | 400M            | 640         | 15         | 32         | 4        | 768          | 8e-5, 1e-4, 3e-4, 5e-4, 8e-4, 1.5e-3       |
| <i>Test Set</i>      |              |                 |             |            |            |          |              |  |
| Qwen3-MoE-4B-A0.5B   | 4B           | 530M            | 768         | 18         | 32         | 4        | 768          | 8e-5, 3e-4, 5e-4, 8e-4, 1.5e-3             |
| Qwen3-MoE-12B-A1.3B  | 12B          | 1.3B            | 1280        | 30         | 32         | 4        | 768          | -  |

Table 2: Overview of Qwen3-MoE Model Architectures and Hyperparameters in Section 3.2

| Models                   | Total Params | Activate Params | Hidden Size | Num Layers | Attn Heads | KV Heads | Interm. Size | Learning Rate                      | std                                     |
|--------------------------|--------------|-----------------|-------------|------------|------------|----------|--------------|------------------------------------|---|
| Qwen3-MoE-2B-A0.3B-proxy | 2B           | 290M            | 512         | 18         | 32         | 4        | 512          | 8e-5, 1e-4, 3e-4, 5e-4, 1e-3, 2e-3 | 0.01, 0.015, 0.02, 0.03, 0.04           |
| Qwen3-MoE-4B-A0.5B       | 4B           | 530M            | 768         | 18         | 32         | 4        | 768          | -                                  | -                                       |
| Qwen3-MoE-2B-A0.3B-proxy | 2B           | 290M            | 640         | 18         | 32         | 4        | 384          | 1e-4, 3e-4, 5e-4, 1e-3, 2e-3       | 0.0005, 0.001, 0.002, 0.005, 0.01, 0.02 |
| Qwen3-MoE-12B-A1.3B      | 12B          | 1.3B            | 1280        | 30         | 32         | 4        | 768          | -                                  | -                                       |

Table 3: Overview of Qwen3-MoE Model Architectures and Hyperparameters in Section 6.1

| Models               | Total Params | Activate Params | Hidden Size | Num Layers | Attn Heads | KV Heads | Interm. Size | Learning Rate                                       |
|----------------------|--------------|-----------------|-------------|------------|------------|----------|--------------|---|
| <i>Training Set</i>  |              |                 |             |            |            |          |              |   |
| Qwen3-MoE-0.5B-A0.1B | 550M         | 100M            | 256         | 3          | 32         | 4        | 768          | 8e-5, 3e-4, 8.75e-4, 1e-3, 1.5e-3, 2e-3, 3e-3, 4e-3 |
| Qwen3-MoE-1B-A0.2B   | 1B           | 190M            | 384         | 9          | 32         | 4        | 768          | 8e-5, 3e-4, 7.24e-4, 1e-3, 1.5e-3, 2e-3, 3e-3, 4e-3 |
| Qwen3-MoE-2B-A0.3B   | 2B           | 280M            | 512         | 12         | 32         | 4        | 768          | 8e-5, 3e-4, 6.36e-4, 1e-3, 1.5e-3, 2e-3, 3e-3, 4e-3 |
| Qwen3-MoE-3B-A0.4B   | 3B           | 400M            | 640         | 15         | 32         | 4        | 768          | 8e-5, 3e-4, 5.90e-4, 1e-3, 1.5e-3, 2e-3, 3e-3, 4e-3 |
| Qwen3-MoE-4B-A0.5B   | 4B           | 530M            | 768         | 18         | 32         | 4        | 768          | 8e-5, 3e-4, 5.55e-4, 1e-3, 1.5e-3, 2e-3, 3e-3, 4e-3 |

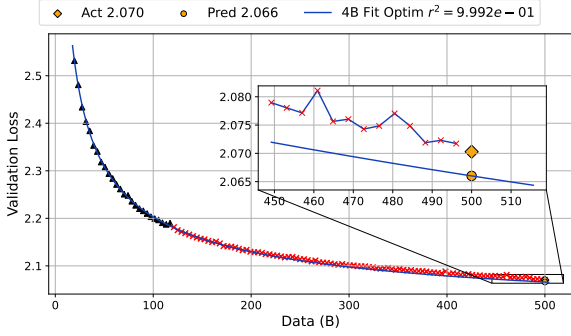


Figure 13: The precision of Equation 4 on 4B Fit Optim model.

879 lation via Equation 4 to be reasonable within the  
880 data range discussed in this paper.

### 881 A.3 Details of Fitting Experiments

#### 882 A.3.1 Global Optimal Learning Rate

883 This subsection contains the whole fitting process  
884 of Section 3.1.

885 The validation loss curves for various models un-  
886 der different learning rates, derived from the global  
887 optimal learning rate search experiments, are il-  
888 lustrated in Figure 14. We utilize the smoothed  
889 data via Equation 4 as the input for subsequent fit-  
890 ting stages. Furthermore, we employ this equation  
891 to extrapolate the validation loss for each model  
892 at a training volume of 200B tokens across dis-  
893 tinct learning rates, thereby augmenting the dataset  
894 available for fitting.

895 We sample validation loss data points at 10B-  
896 token intervals, ranging from 80B to 220B tokens,  
897 to facilitate subsequent analysis and curve fitting.  
898 Figure 15 illustrates the variation of validation loss  
899 as a function of the learning rate  $\eta$  with different  
900 model size and data size. Figure 16 presents a 3-D  
901 visualization of the relationship among loss, learn-  
902 ing rate, and training data size. Upon observing a  
903 distinct local minimum, we employ the quadratic  
904 polynomial defined in Equation 2 to fit the data  
905 (as shown in Figure 17). The coefficients of deter-  
906 mination ( $R^2$ ) consistently exceed 0.995, enabling  
907 a precise estimation of the optimal learning rate  
908 based on these curves.

909 As shown in Figure 18(a) and 18(b), the global  
910 optimal learning rate exhibits a power-law relation-  
911 ship with both the model parameter count  $N$  and  
912 training data size  $D$ . With reference to the stud-  
913 ies of Bjorck et al. (2025), we decide to use the  
914 following functional form:

$$\eta_{opt}(N, D) = C_{\eta} \cdot N^{-\alpha} \cdot D^{-\beta}, \quad (8) \quad 915$$

916 where  $C_{\eta}, \alpha, \beta$  are positive constants. After em-  
917 ploying non-linear least squares to fit the curve, we  
918 finally get the parameters of Equation 6:

$$C_{\eta} \sim 38.4588, \alpha \sim 0.2219, \beta \sim 0.3509. \quad (9) \quad 919$$

#### 920 A.3.2 Module-Level Optimal Learning Rates

921 This subsection details the step-by-step process of  
922 searching Module-Level Optimal LR.

923 We split the model into the following four  
924 groups of parameters:

- 925 • **Embedding Parameters**, which is the word  
926 embedding layer of a model,
- 927 • **Hidden Parameters**, mainly composed of  
928 self-attention and layer norm modules,
- 929 • **Router**, which contains the router matrix and  
930 experts,
- 931 • **LM Head Parameters**, which is the unem-  
932 bedding output layer.

933 Similar to our experiments in Section A.3.2,  
934 while searching optimal LR across different mod-  
935 ule groups, the training data size in set to approx-  
936 imately 120B tokens. According to the results  
937 above, we can derive the global optimal LR  $\eta_{opt}$   
938 of every size of model in the experiment via Equation  
939 2:

940 In the following stages, we sequentially conduct-  
941 ing experiments in the order of LM Head, Router,  
942 Hidden, and Embedding parameters with greedy  
943 search strategy.

944 **LM Head.** First, we begin with the LM Head  
945 module. By varying the learning rate  $\eta^{out}$  of the  
946 LM Head weights within a specified range while  
947 fixing the learning rates of all other weights to  
948 the current model’s global optimal learning rate(i.e.  
949  $\eta^{emb} = \eta^{hidden} = \eta^{router} = \eta_{opt}$ ), we conduct the  
950 search following the method described in Section  
951 3.1. The curve fitted using Equation 2 is shown in  
952 Figure 20, where the fitted minimum is taken as the  
953 module-level learning rate  $\eta_{opt}^{out}$  for LM Head(Table  
954 5).

955 **Router.** In the second searching stage, we set  
956  $\eta^{emb} = \eta^{hidden} = \eta_{opt}, \eta^{out} = \eta_{opt}^{out}$  and search  
957 learning rate on Router layers. The results are  
958 illustrated in Figure 21 and Table 6

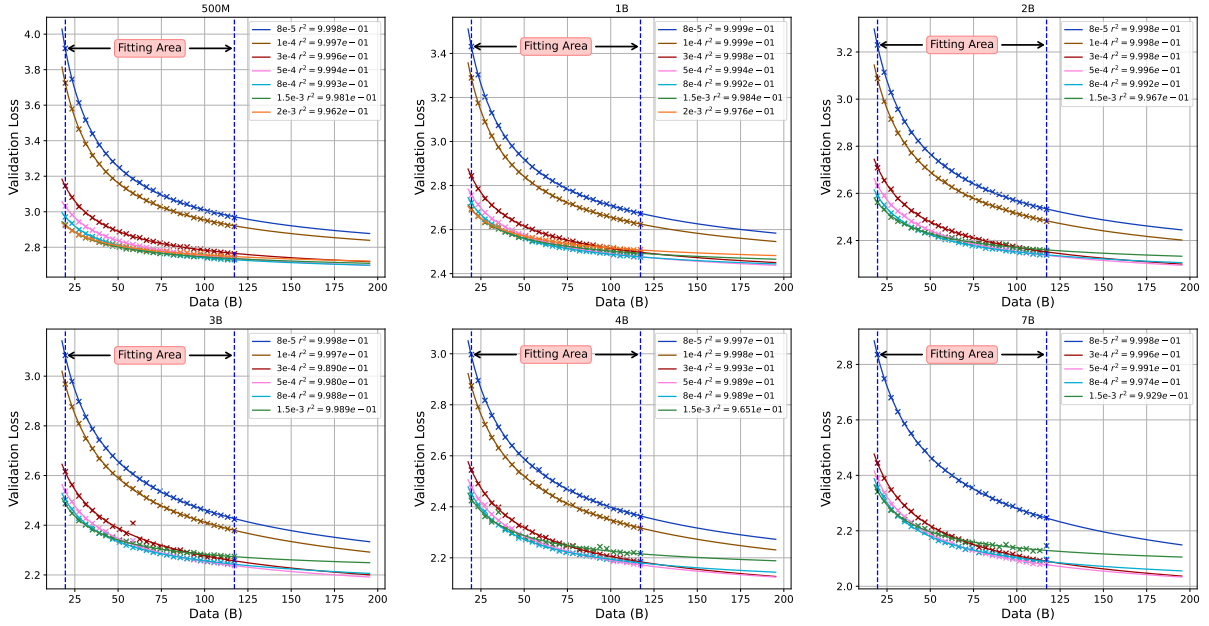


Figure 14: Results of fitting via  $L(D) = L_0 + A \cdot D^{-\gamma}$  for each group of experiments.

Table 4: Global Optimal LR at 120B

| $N$          | 500M    | 1B      | 2B      | 3B      | 4B      |
|--------------|---------|---------|---------|---------|---------|
| $\eta_{opt}$ | 8.75e-4 | 7.24e-4 | 6.36e-4 | 5.90e-4 | 5.55e-4 |

**Hidden.** Next, set  $\eta^{emb} = \eta_{opt}, \eta^{router} = \eta_{opt}^{router}, \eta^{out} = \eta_{opt}^{out}$  while searching optimal learning rate on Hidden parameters to obtain  $\eta_{opt}^{hidden}$ . The results are illustrated in Figure 22 and Table 7

**Embedding.** Finally, we set  $\eta^{router} = \eta_{opt}^{router}, \eta^{hidden} = \eta_{opt}^{hidden}, \eta^{out} = \eta_{opt}^{out}$  and conduct learning rate searching on Embedding layer and get its optimal learning rate  $\eta_{opt}^{emb}$ . The results are illustrated in Figure 23 and Table 8

**Overall Results.** The overall results of module-level optimal learning rate are shown in Table 9

### A.3.3 $\mu$ Transfer

We refer to Mlodozeniec et al. (2025) to conduct our  $\mu$ Transfer experiments. The transfer method is shown in Table 10. As we maintain an invariant batch size across all experimental configurations, we only consider the influence of training token counts alongside model width and depth when employ  $\mu$ Transfer.

Table 5: LM Head Optim LR at 120B

| $N$                | <b>500M</b> | <b>1B</b> | <b>2B</b> | <b>3B</b> | <b>4B</b> |
|--------------------|-------------|-----------|-----------|-----------|-----------|
| $\eta_{opt}^{out}$ | 6.92e-4     | 425e-4    | 3.72e-4   | 3.81e-4   | 2.86e-4   |

Table 6: Router Optim LR at 120B

| $N$                   | <b>500M</b> | <b>1B</b> | <b>2B</b> | <b>3B</b> | <b>4B</b> |
|-----------------------|-------------|-----------|-----------|-----------|-----------|
| $\eta_{opt}^{router}$ | 7.62e-4     | 6.55e-4   | 5.93e-4   | 5.23e-4   | 4.89e-4   |

Table 7: Hidden Optim LR at 120B

| $N$                   | <b>500M</b> | <b>1B</b> | <b>2B</b> | <b>3B</b> | <b>4B</b> |
|-----------------------|-------------|-----------|-----------|-----------|-----------|
| $\eta_{opt}^{hidden}$ | 1.05e-3     | 6.99e-4   | 5.80e-4   | 4.79e-4   | 3.78e-4   |

Table 8: Embedding Optim LR at 120B

| $N$                | <b>500M</b> | <b>1B</b> | <b>2B</b> | <b>3B</b> | <b>4B</b> |
|--------------------|-------------|-----------|-----------|-----------|-----------|
| $\eta_{opt}^{out}$ | 1.95e-3     | 8.38e-4   | 2.00e-3   | 1.15e-3   | 1.05e-3   |

Table 9: Module-Level Optim LR at 120B

| $N$                   | <b>500M</b> | <b>1B</b> | <b>2B</b> | <b>3B</b> | <b>4B</b> |
|-----------------------|-------------|-----------|-----------|-----------|-----------|
| $\eta_{opt}^{out}$    | 6.92e-4     | 425e-4    | 3.72e-4   | 3.81e-4   | 2.86e-4   |
| $\eta_{opt}^{router}$ | 7.62e-4     | 6.55e-4   | 5.93e-4   | 5.23e-4   | 4.89e-4   |
| $\eta_{opt}^{hidden}$ | 1.05e-3     | 6.99e-4   | 5.80e-4   | 4.79e-4   | 3.78e-4   |
| $\eta_{opt}^{out}$    | 1.95e-3     | 8.38e-4   | 2.00e-3   | 1.15e-3   | 1.05e-3   |

Table 10: Hyperparameters' transfer rule of  $\mu$ Transfer

|                  | <b>Parameterisation:</b>    |              | $\mu\mathbf{P}$                            | <b>Complete<sup>(d)</sup>P</b>                           |                               |
|------------------|-----------------------------|--------------|--|--|-------------------------------|
| Multipliers      | MHA Residual                |              | $\mathbf{x} + \text{MHABlock}(\mathbf{x})$ | $\mathbf{x} + m_L^{-\alpha} \text{MHABlock}(\mathbf{x})$ |                               |
|                  | MLP Residual                |              | $\mathbf{x} + \text{MLPBlock}(\mathbf{x})$ | $\mathbf{x} + m_L^{-\alpha} \text{MLPBlock}(\mathbf{x})$ |                               |
|                  | Unemb. Fwd                  |              | Unaugmented                                | Unaugmented  |                               |
| Init Variances   | Input Emb.                  |              |  |  |                               |
|                  | Hidden weights              |              | $\times m_N^{-1}$                          | $\times m_N^{-1}$  |                               |
|                  | Hidden biases/norms         | $\sigma_b^2$ |  |  |                               |
|                  | Unemb. LN                   |              |  |  |                               |
| Learning Rates   | Unemb. Weights              |              | $\times m_N^{-2}$                          | $\times m_N^{-2}$  |                               |
|                  | Input Emb.                  |              |  |  |                               |
|                  | Hidden weights              |              | $\times m_N^{-1}$                          | $\times m_N^{-1} \times m_L^{\alpha-1}$                  | $\times \sqrt{\frac{1}{m_D}}$ |
|                  | Hidden biases/norm          | $\eta_b$     |  | $\times m_L^{\alpha-1}$                                  |                               |
| Unemb. LN        |                             |              |  |  |                               |
| Unemb. weights   |                             |              | $\times m_N^{-1}$                          | $\times m_N^{-1}$  |                               |
| AdamW $\epsilon$ | Hidden weights/biases/norms |              | $\times m_N^{-1}$                          | $\times m_N^{-1} \times m_L^{-\alpha}$                   | $\times \sqrt{m_D}$           |
|                  | QK norms                    |              | NA   | $\times m_L^{-\alpha}$                                   |                               |
|                  | Input Emb.                  | $\epsilon_b$ | $\times m_N^{-1}$                          | $\times m_N^{-1}$  |                               |
|                  | Output weights/biases/norms |              |  |  |                               |
| Weight decay     | Hidden weights              |              | $\times m_N$                               | $\times m_N$   | $\times \sqrt{\frac{1}{m_D}}$ |
|                  | Unemb. weights              | $\lambda_b$  | $\times m_N$                               | $\times m_N$   |                               |
|                  | Rest                        |              | $\times 1$                                 | $\times 1$   |                               |

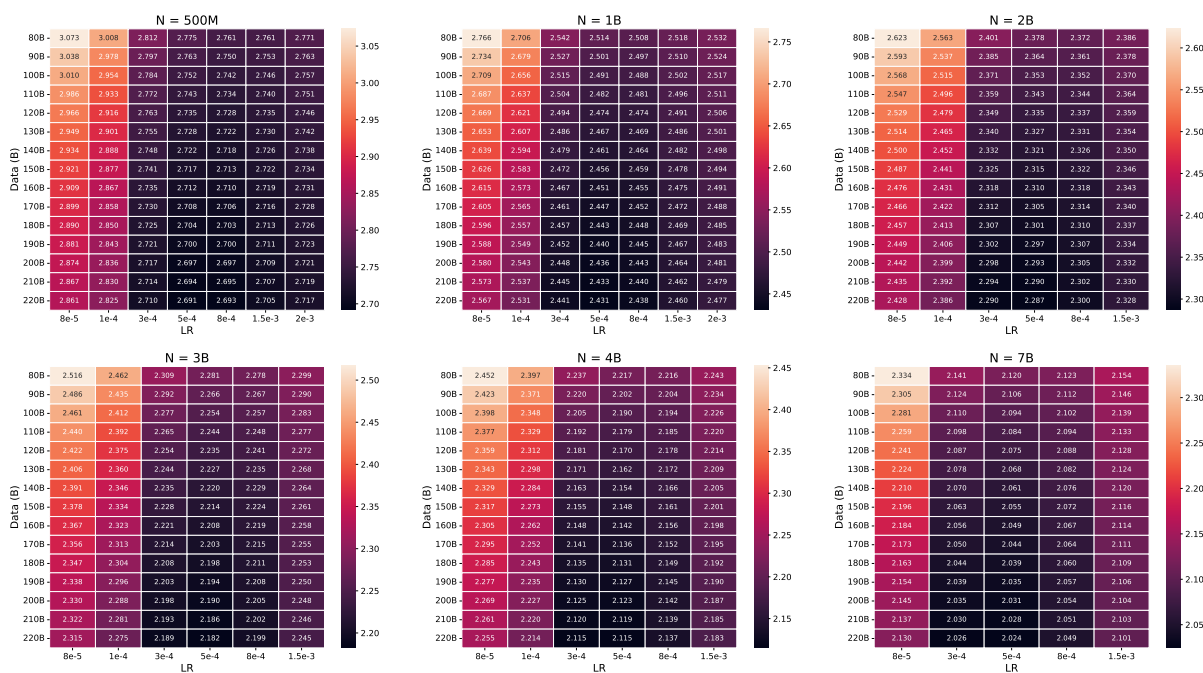


Figure 15: Relationship among loss, learning rate, and training data size of various model.

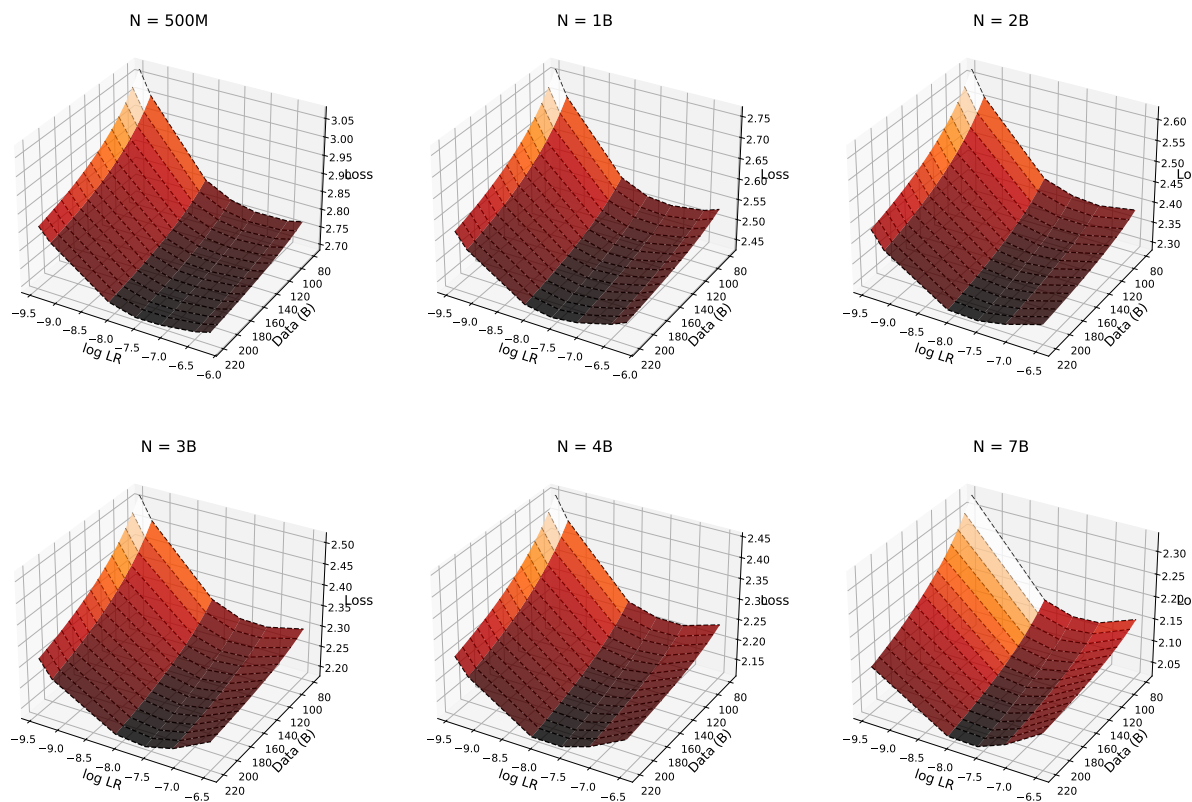


Figure 16: 3D visualization of 15

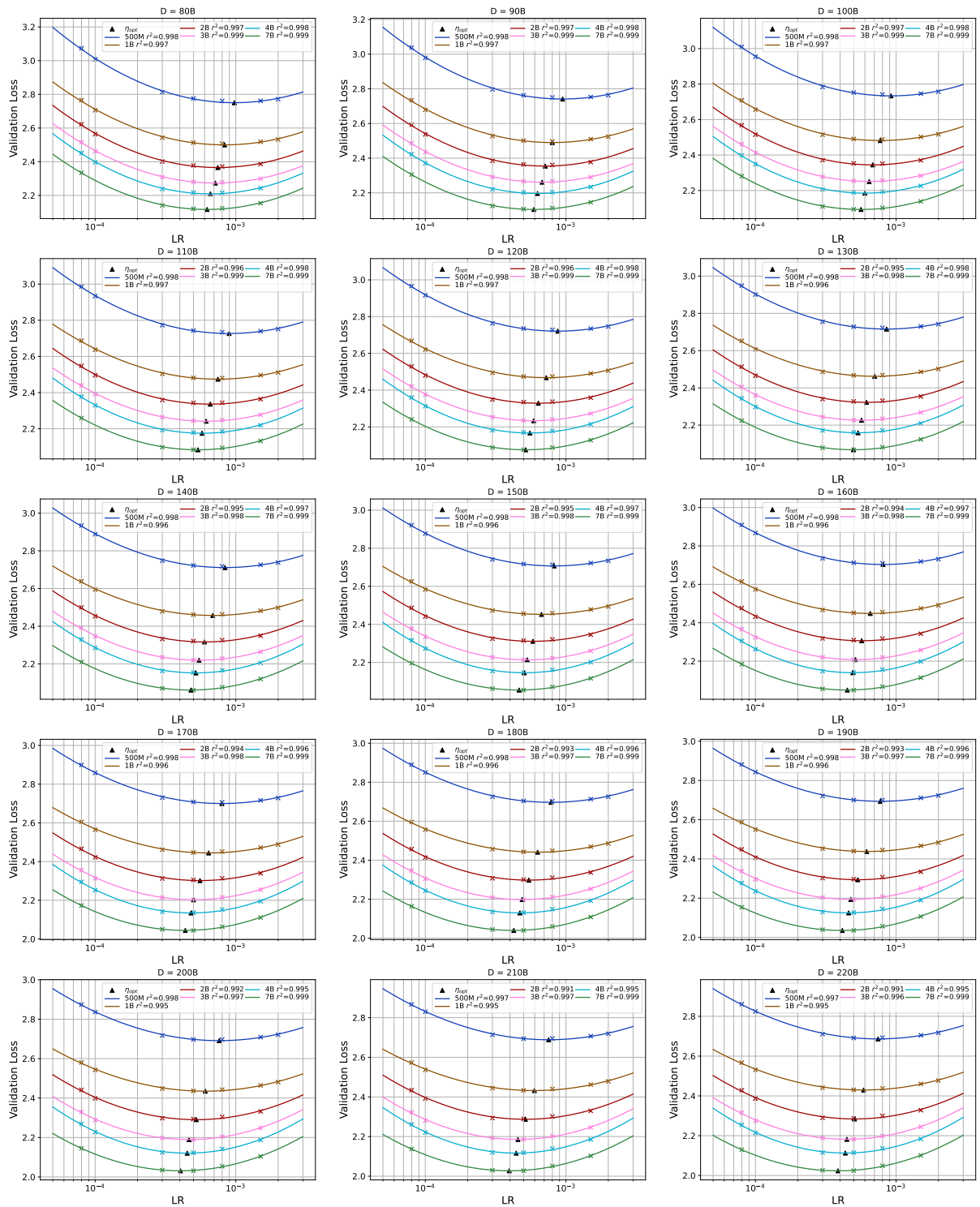


Figure 17: All results of fitting with Equation 2 across different amount of data.

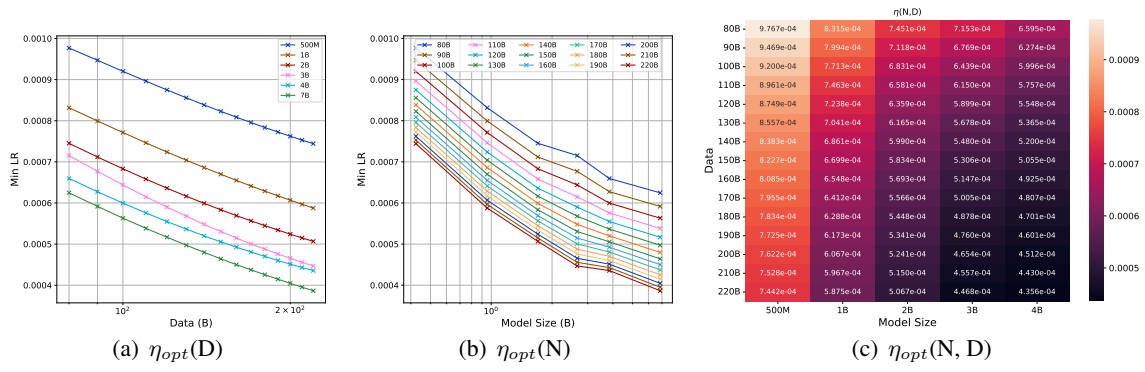


Figure 18: Relationship among learning rate  $\eta$ , model size  $N$  and training data size  $D$ .

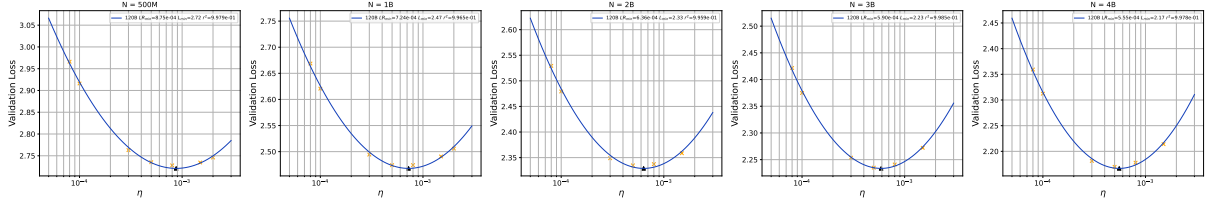


Figure 19: Initial stage of module-level learning rate searching.

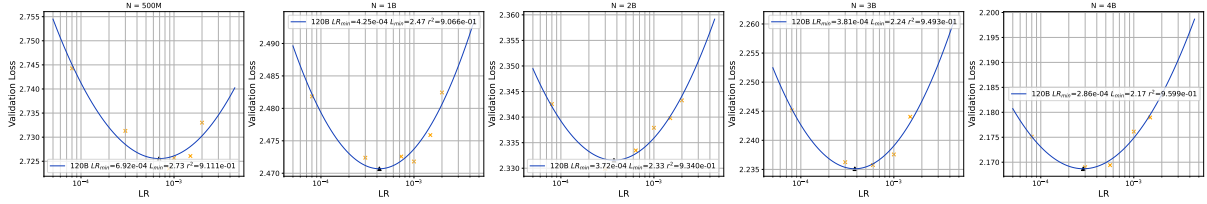


Figure 20: LM Head stage of module-level learning rate searching.

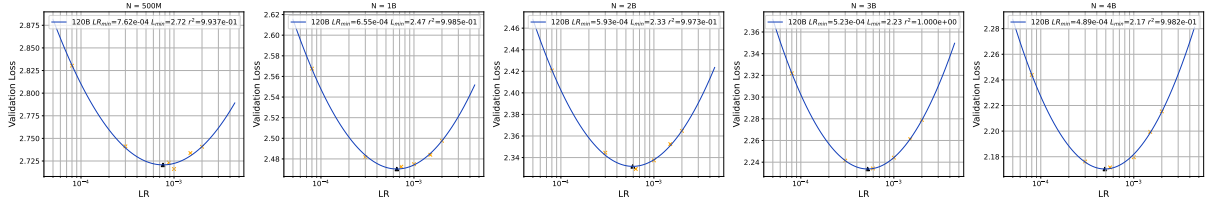


Figure 21: Router stage of module-level learning rate searching.

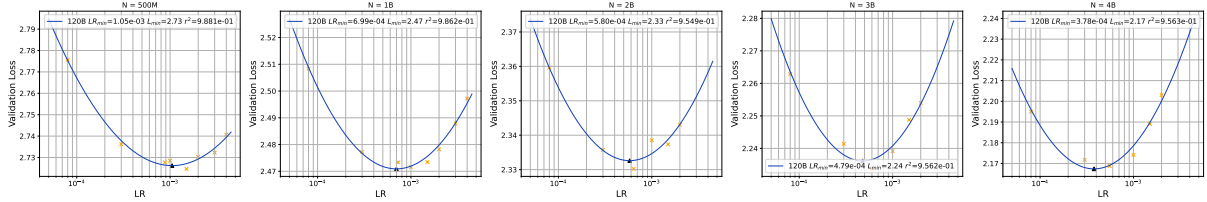


Figure 22: Hidden stage of module-level learning rate searching.

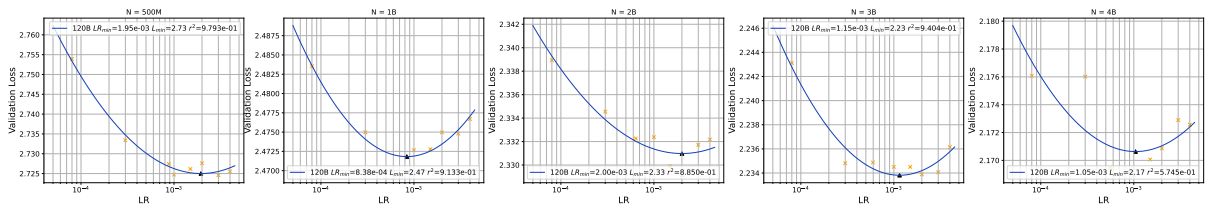


Figure 23: Embedding stage of module-level learning rate searching.