

SPARKLING: Balancing Signal Preservation and Symmetry Breaking for Width-Progressive Learning

Qifan Yu^{1,2,*}, Xinyu Ma¹, Zhijian Zhuo¹, Minrui Wang¹, Deyi Liu¹, Shiyi Zhan¹,
Yiyuan Ma¹, Liang Xiang¹, Xingyan Bin^{1,†}, Di He^{2,†}

¹ByteDance Seed, ²Peking University

*Work done at ByteDance Seed, [†]Corresponding authors

Abstract

Progressive Learning (PL) reduces pre-training computational overhead by gradually increasing model scale. While prior work has extensively explored depth expansion, width expansion remains significantly understudied, with the few existing methods limited to the early stages of training. However, expanding width during the mid-stage is essential for maximizing computational savings, yet it remains a formidable challenge due to severe training instabilities. Empirically, we show that naive initialization at this stage disrupts activation statistics, triggering loss spikes, while copy-based initialization introduces gradient symmetry that hinders feature diversity. To address these issues, we propose **SPARKLING** (balancing **S**ignal **P**reservation **A**nd **s**ymmet**R**y **b**rea**K**ing for width-progressive **L**earn**I**NG), a novel framework for mid-stage width expansion. Our method achieves signal preservation via RMS-scale consistency, stabilizing activation statistics during expansion. Symmetry breaking is ensured through asymmetric optimizer state resetting and learning rate re-warmup. Extensive experiments on Mixture-of-Experts (MoE) models demonstrate that, across multiple width axes and optimizer families, SPARKLING consistently outperforms training from scratch and reduces training cost by up to 35% under 2× width expansion.

Date: February 3, 2026

Correspondence: Qifan Yu at qifanyu@stu.pku.edu.cn, Xingyan Bin at binxingyan@bytedance.com,
Di He at di_he@pku.edu.cn

1 Introduction

Training Large Language Models (LLMs) remains prohibitively expensive, motivating a growing line of work on Progressive Learning (PL) [10, 17, 33], which aims to expand the parameter scale gradually during training instead of training the full scale from scratch [12]. Existing PL methods have demonstrated notable success in both saving training computation and performance enhancement, especially through depth-oriented strategies such as layer stacking [10, 12, 17], block insertion [33, 39], or gradual network growth [37].

Width is another crucial dimension for scaling model parameters [16]. Existing studies have made only limited progress in this direction, and a general and systematic mechanism has yet to be established [5, 6, 40–42, 45]. More importantly, previous investigations have been largely limited to expansion during the initial portion of training, e.g., less than 10–30% of training tokens [10, 28]. Such early expansion offers negligible computational advantages over training the target-width model from scratch and fundamentally undermines the primary motivation of PL—reducing training costs. To make PL practically viable, width expansion must be conducted

during the intermediate stage of pre-training (i.e., mid-stage expansion). However, it is precisely in this regime that width expansion becomes challenging. We hereby analyze the core challenges as follows.

We identify the first core challenge as *signal preservation* during mid-stage expansion. Prior width-based PL work has largely treated “preservation” as *loss continuity* at the expansion point, i.e., function preservation (FP) where the expanded model is initialized to match the pre-expansion mapping [5, 6, 13, 28, 31]. While FP is useful, we argue that the core mechanism behind the scenes is whether the expansion preserves the statistical distribution of intermediate activations, most notably the *root-mean-square (RMS)* scale of hidden representations [44]. Concretely, RMS-scale mismatch alters layerwise signal magnitudes and propagates through residual streams; this destabilizes optimization even when no instantaneous loss spike occurs at the moment of expansion [1, 38].

Based on this perspective, we design RMS-preserving strategies for several standard initialization regimes—copy, random, and zero—ensuring each can be applied without inducing activation-scale shocks. Interestingly, these RMS-preserved variants reveal a counter-intuitive limitation of copying: despite its appeal for function preservation, copy-based expansion can underperform compared to RMS-preserved random or zero initialization in terms of post-expansion recovery. This observation indicates that, beyond maintaining loss and activation scale, some additional and critical issues for copy-based expansion remain unresolved.

This brings us to the second core challenge: copy-based expansion, while strongest in forward continuity, induces *backward symmetry* [6]. Duplicating channels creates duplicated parameter subspaces that receive identical gradients and thus evolve identically, leaving the new capacity functionally redundant [34]. Crucially, this symmetry is not an artifact of a specific optimizer: it arises under both element-wise optimizers such as AdamW [20] as well as spectral-style updates such as Muon [15, 19], causing a persistent coupled state in which copied components fail to diversify.

Motivated by these observations, we frame mid-stage width expansion as balancing two complementary principles: *Signal Preservation* and *Symmetry Breaking*. On the preservation side, we enforce RMS-scale consistency during expansion, ensuring that the expanded model maintains stable hidden-state statistics and thereby supports smooth post-expansion optimization. On the symmetry side, we introduce targeted interventions that act only in the backward dynamics while leaving the copied forward function intact: (i) a controlled optimizer-state reset for newly introduced parameters to remove inherited symmetric momentum, and (ii) an asymmetric learning rate re-warmup schedule that selectively stimulates the new parameters without globally perturbing the well-adapted pre-trained ones. Together, these mechanisms preserve forward continuity at the moment of expansion, while inducing sufficient asymmetry in subsequent updates for the expanded capacity to diverge and encode meaningful features.

In summary, our contributions can be highlighted as follows:

- We investigate the challenges of width expansion during the critical **mid-stage of pre-training**, a regime largely unexplored in prior work due to stability concerns. We identify that successful expansion hinges on two complementary principles: *Signal Preservation* to stabilize activation statistics, and *Symmetry Breaking* to resolve gradient coupling in copy-based initialization.
- We propose SPARKLING, a practical framework that implements both principles through a suite of concrete mechanisms—including RMS-scale consistency, copy-based initialization, asymmetric optimizer state resetting, and asymmetric learning rate re-warmup—that jointly resolve the optimization challenges inherent to expanding deep within the pre-training trajectory.
- We empirically validate the generality of SPARKLING across multiple width axes (including hidden dimension and MoE expert intermediate dimension) and optimizer families (including AdamW and Muon). Under a fixed token budget, our PL approach consistently **outperforms training the full-scale model from scratch** on downstream evaluations, while **reducing training costs by up to 35 %** when scaling to $2\times$ width, demonstrating both **effectiveness** and **efficiency** of SPARKLING.

2 Related Work

Progressive learning (PL) has emerged as a resource-efficient paradigm that accelerates training by gradually expanding the architecture from a small base model to a target scale during training [5, 6, 12, 17]. From the perspective of depth expansion, existing strategies typically grow by stacking layers [10, 12, 17] or inserting blocks [33, 39]. Existing approaches for width expansion largely prioritize function preservation (FP) via parameter mapping [6], advanced initialization schemes like AKI and its variants [5, 40, 45], or temporarily masking new structures [41]. To address redundancy from simple copying [6], various heuristic interventions have been adopted, including uneven splitting [6, 10, 31] and symmetric perturbations [34, 35, 42].

Beyond initialization, significant effort has been directed toward stabilizing post-growth optimization dynamics: Wang et al. [31] advocate for accelerated decay schedules on the premise that expanded models start closer to local optima, while Yuan et al. [42] utilize weight-norm to rebalance gradient contributions and Shen et al. [28] propose dynamics-preserving growth operators to align the expanded model’s loss trajectory. Other methods attempt to learn growth operators [24, 30] or construct gradient-maximizing weights [11].

However, these strategies typically address either forward initialization or backward optimization dynamics in isolation. In this work, we innovatively establish a systematic framework balancing both perspectives. We first argue that the mechanism underlying widely used *function-preserving* initializations is fundamentally *RMS preservation*, and then redesign the optimization procedures to address the symmetry issues that inevitably arise from such preservation-focused initialization strategies.

3 RMS Scale Consistency of Activation

3.1 Why RMS Mismatch Destabilizes Training

We start by defining the root-mean-square (RMS) magnitude of a vector $\mathbf{h} \in \mathbb{R}^d$ as [44]

$$\text{RMS}(\mathbf{h}) := \frac{\|\mathbf{h}\|_2}{\sqrt{d}} = \sqrt{\frac{1}{d} \sum_{i=1}^d h_i^2}. \quad (1)$$

Consider a linear layer

$$\mathbf{y} = \mathbf{W}\mathbf{x}, \quad \mathbf{W} \in \mathbb{R}^{d_{\text{out}} \times d_{\text{in}}}, \quad \mathbf{x} \in \mathbb{R}^{d_{\text{in}}}, \quad \mathbf{y} \in \mathbb{R}^{d_{\text{out}}}. \quad (2)$$

where \mathbf{x} and \mathbf{y} denote the input and output hidden states, respectively. Our focus is the RMS scale of the activations:

$$r := \frac{s_{\text{out}}}{s_{\text{in}}} = \frac{\text{RMS}(\mathbf{y})}{\text{RMS}(\mathbf{x})}, \quad r^{(\text{post})} = r^{(\text{pre})}, \quad (3)$$

requiring this scale to remain unchanged after expansion.

We enforce the RMS invariance in Eq. (3) as a signal preservation constraint. A trained Transformer block implicitly defines an operating regime via its input-output statistics, within which representations are well-formed and features remain meaningful. If width expansion perturbs activation RMS during expansion, post-expansion hidden states can drift away from the pre-expansion scale manifold, causing subsequent blocks to receive out-of-regime inputs. RMS-preserving expansion mitigates this shift by keeping block-wise input/output magnitudes within the original domain, thereby maintaining the fidelity and generalization of the pre-trained function immediately after expansion.

In modern LLMs that adopt pre-normalization (e.g., Qwen3 [36], DeepSeek-V3 [9], OLMoE [23]), RMS-preserving becomes even more critical because the update explicitly couples the residual stream with the branch output. Concretely, in a typical residual block with pre-norm, the hidden state is updated as

$$\mathbf{h} \leftarrow \mathbf{h} + f(\text{Norm}(\mathbf{h})), \quad (4)$$

where $f(\cdot)$ denotes a residual branch such as an attention or MLP sublayer, and $\text{Norm}(\cdot)$ refers to token-wise feature normalization, i.e., LayerNorm or its variants used in LLMs, most notably RMSNorm [44].

Here, pre-norm stabilizes the input to $f(\cdot)$ but does not constrain its output scale, so the residual dynamics depend on the ratio $\text{RMS}(f(\text{Norm}(\mathbf{h})))/\text{RMS}(\mathbf{h})$. After expansion, an RMS mismatch shifts the calibrated mixing between the main path and the transformed branch, making it either overwhelming on the main stream or nearly identity. By preserving RMS through expansion, we keep layerwise dynamics coherent and maintain the balanced residual regime.

3.2 RMS-Preserving Expansion

We continue to discuss RMS-preserving width growth in three cases: (i) *fan-out* expansion, which expands the output dimension (d_{out}), (ii) *fan-in* expansion, which expands the input dimension (d_{in}), and (iii) *RMSNorm weight* expansion, which widens the RMSNorm scale.

In practice, fan-out and fan-in expansions typically appear as a paired transformation across two consecutive layers that share an intermediate width. For example, in an MLP that widens the expert intermediate dimension, the *up* and *gate* projections become fan-out expansions¹, and the subsequent *down* projection becomes fan-in expansion. Similarly, in attention, the *vhead* projection and the output projection likewise form such a paired transformation. Therefore, whenever we widen any width dimension, the two sides are naturally correlated and should be discussed jointly in pairs.

3.2.1 Preliminaries

Under the linear layer defined in Eq. (2), the output activation RMS is given by

$$\text{RMS}(\mathbf{y}) = \sqrt{\frac{1}{d_{\text{out}}} \|\mathbf{y}\|_2^2} = \sqrt{\frac{1}{d_{\text{out}}} \sum_{i=1}^{d_{\text{out}}} y_i^2}. \quad (5)$$

Leveraging the property of high-dimensional isotropy in wide neural networks, where feature vectors tend toward asymptotic orthogonality [2, 27], we can assume that $\{y_i\}_{i=1}^{d_{\text{out}}}$ are identically distributed and satisfy $\mathbb{E}[y_i] = 0$. Taking expectation over the data yields

$$\mathbb{E}[\text{RMS}^2(\mathbf{y})] = \mathbb{E}\left[\frac{1}{d_{\text{out}}} \sum_{i=1}^{d_{\text{out}}} y_i^2\right] = \mathbb{E}[y_i^2] = \text{Var}(y_i). \quad (6)$$

Therefore, when the input RMS s_{in} is kept unchanged, preserving the output RMS scale is equivalent to preserving the per-coordinate variance

$$\text{Var}(y_i) = d_{\text{in}} \sigma_w^2 \sigma_x^2, \quad (7)$$

where σ_w^2 and σ_x^2 denote the shared variances of w_{ij} and x_j , respectively. We derive Eq. (7) in Appendix A.1.

3.2.2 Fan-out Expansion

In fan-out expansion, the output dimension grows from d_{out} to d'_{out} ($> d_{\text{out}}$) while keeping d_{in} unchanged, denoted by

$$\mathbf{y}' = \mathbf{W}' \mathbf{x}, \quad \mathbf{W}' = \begin{bmatrix} \mathbf{W} \\ \tilde{\mathbf{W}} \end{bmatrix} \in \mathbb{R}^{d'_{\text{out}} \times d_{\text{in}}}, \quad (8)$$

$$\mathbf{y}' = \begin{bmatrix} \mathbf{y} \\ \tilde{\mathbf{y}} \end{bmatrix} \in \mathbb{R}^{d'_{\text{out}}}, \quad \tilde{\mathbf{y}} = \tilde{\mathbf{W}} \mathbf{x}. \quad (9)$$

Naturally, fan-out expansion preserves activation RMS as long as the newly introduced output channels \mathbf{W}' are distributionally consistent with the pre-expansion ones. Concretely, when the added rows are initialized by *copying* or by *randomly sampling* from the same distribution as the original weights, the expanded output \mathbf{y}' remains distributionally aligned with \mathbf{y} and thus retains the same RMS scale.

¹We thus treat the *gate* projection in the same way as the *up* projection under RMS-preserving expansion, and regard the resulting gate activation output as $\Theta(1)$ multiplicative factor in expectation.

3.2.3 Fan-in Expansion

In fan-in expansion, the input dimension grows from d_{in} to d'_{in} ($> d_{\text{in}}$) while keeping d_{out} unchanged, denoted by

$$\mathbf{y}' = \mathbf{W}' \mathbf{x}', \quad \mathbf{W}' = \alpha [\mathbf{W} \quad \tilde{\mathbf{W}}] \in \mathbb{R}^{d_{\text{out}} \times d'_{\text{in}}}, \quad (10)$$

$$\mathbf{x}' = \begin{bmatrix} \mathbf{x} \\ \tilde{\mathbf{x}} \end{bmatrix} \in \mathbb{R}^{d'_{\text{in}}}, \quad \mathbf{y}' = \alpha(\mathbf{W}\mathbf{x} + \tilde{\mathbf{W}}\tilde{\mathbf{x}}). \quad (11)$$

RMS-preserving fan-in expansion seeks the scaling factor α satisfying the invariance of Eq. (7) across expansion.

Random or One-Side Copied. If the newly added fan-in coordinates are initialized by same-distribution random sampling, or by copying on only one side (i.e., only \mathbf{W} or only \mathbf{x} is copied while the other remains random), then Eq. (6) and its underlying assumptions continue to hold, resulting in a shared per-coordinate variance after expansion:

$$\text{Var}(y'_i) = \sum_{j=1}^{d'_{\text{in}}} \text{Var}(w'_{ij} x'_j) = \sum_{j=1}^{d'_{\text{in}}} \sigma_{w'}^2 \sigma_{x'}^2 = d'_{\text{in}} \sigma_{w'}^2 \sigma_{x'}^2. \quad (12)$$

With unchanged input scale $\sigma_{x'}^2 = \sigma_x^2$, variance preservation requires

$$d'_{\text{in}} \sigma_{w'}^2 = d_{\text{in}} \sigma_w^2 \implies \sigma_{w'} = \sqrt{\frac{d_{\text{in}}}{d'_{\text{in}}}} \sigma_w, \quad (13)$$

which implies that the weights should be rescaled as

$$w'_{ij} = \sqrt{\frac{d_{\text{in}}}{d'_{\text{in}}}} w_{ij}, \quad \forall i = 1, \dots, d_{\text{out}}, \quad j = 1, \dots, d'_{\text{in}}, \quad (14)$$

thereby keeping $\text{Var}(y'_i) = \text{Var}(y_i)$ and the output RMS invariant under fan-in expansion.

Both-Sides Copied. A qualitatively different regime arises when *both* sides of the newly introduced fan-in coordinates are created by copying existing dimensions, where the independence across fan-in dimensions is violated.

Let c denote the copy ratio. $0 < c \leq 1$ corresponds to the setting where each copied dimension is duplicated exactly once, while $c > 1$ corresponds to that some dimensions may be copied multiple times. Generally, we have

$$d'_{\text{in}} = (1 + c)d_{\text{in}}. \quad (15)$$

The invariance of Eq. (7) requires the weights be rescaled as

$$w'_{ij} = \begin{cases} \frac{1}{\sqrt{1+3c}} w_{ij}, & 0 < c \leq 1, \\ \frac{1}{1+c} w_{ij}, & c > 1, \end{cases} \quad (16)$$

or equivalently,

$$w'_{ij} = \begin{cases} \sqrt{\frac{d_{\text{in}}}{3d'_{\text{in}} - 2d_{\text{in}}}} w_{ij}, & d_{\text{in}} < d'_{\text{in}} \leq 2d_{\text{in}}, \\ \frac{d_{\text{in}}}{d'_{\text{in}}} w_{ij}, & d'_{\text{in}} > 2d_{\text{in}}, \end{cases} \quad (17)$$

$\forall i = 1, \dots, d_{\text{out}}, \quad j = 1, \dots, d'_{\text{in}}$. We provide the full derivation of the above scaling factor in Appendix A.2.

One-Side Zero. Empirically, we find that RMS-preserving expansion should treat the zero-initialized side as *random* rather than strictly loss preserving at the expansion moment, and we include detailed analysis in Appendix C.

3.2.4 RMSNorm Weight Expansion

We next discuss how to expand the RMSNorm scale when widening the dimension. For RMSNorm parameterized by $\gamma \in \mathbb{R}^d$, omitting the ϵ term for clarity, we have

$$\mathbf{z} = \text{RMSNorm}(\mathbf{x}; \gamma) = \frac{\mathbf{x} \odot \gamma}{\text{RMS}(\mathbf{x})}, \quad z_i = \frac{x_i \gamma_i}{\text{RMS}(\mathbf{x})}. \quad (18)$$

Applying Eq. (6) to \mathbf{z} yields

$$\mathbb{E}[\text{RMS}^2(\mathbf{z})] = \text{Var}(z_i) = \frac{1}{\text{RMS}^2(\mathbf{x})} \sigma_x^2 \sigma_\gamma^2 \sim \sigma_\gamma^2, \quad (19)$$

where $\sigma_x^2 := \text{Var}(x_i)$ and $\sigma_\gamma^2 := \text{Var}(\gamma_i)$, and Eq. (6) along with its underlying assumptions is used for both \mathbf{x} and \mathbf{z} .

Therefore, preserving the output RMS of \mathbf{z} under width expansion is effectively equivalent to preserving the RMS of parameter γ . Thus, when expanding RMSNorm from d to $d' > d$, initializing the new coordinates of γ by copying or randomly sampling from the same distribution naturally maintains $\text{RMS}(\mathbf{z})$ without any additional rescaling.

3.3 RMS-Preserving Expansion Improves Late-Stage Convergence.

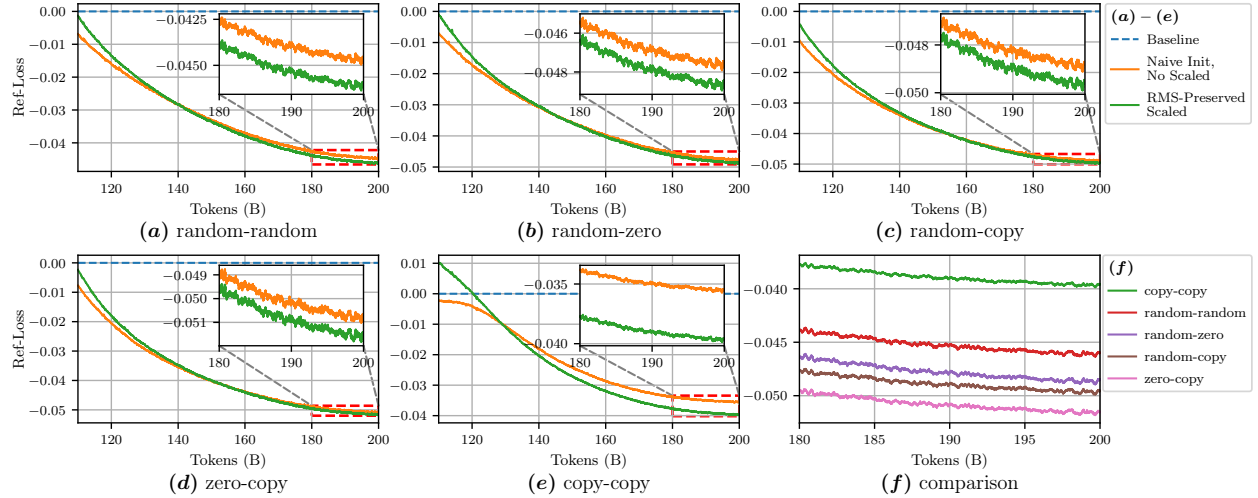


Figure 1 RMS-preserving rescaling consistently improves late-stage convergence under MoE expert inner-dimension expansion. We expand the expert inner dimension from 512 \rightarrow 1024 at 100B tokens and plot *reference-loss* (relative to the pre-expansion reference) over the remaining training tokens. **(a)–(e)** sweep five (*up_proj* – *down_proj init*) pairs. In every case, *Naive Init, No Scaled* yields a smaller immediate loss gap, while *RMS-Preserved Scaled* overtakes later and converges to a lower final loss. **(f)** compares the RMS-preserved late-stage results and highlights a notable pattern: *both-sides copied* significantly underperforms other RMS-preserved strategies.

Experimental Setup. We conduct progressive-learning experiments on OLMoE [22] with 0.5 B active parameters and 2.5 B total parameters, trained for 200 B tokens in total using AdamW optimizer. We perform a mid-stage width expansion at 100 B tokens and then continue to train the expanded model for the remaining 100 B tokens under the same training recipe. Details of the experimental setup are provided in Appendix B.

We consider two width-growth axes: (i) *Inner* $2\times$, which doubles the MoE expert intermediate dimension from 512 to 1024, and (ii) *Hidden* $2\times$, which doubles the model hidden size from 1024 to 2048². In each setting, we compare two initialization strategies for the newly introduced channels: (1) *Naive Init, No Scaled*,

²We decouple hidden dimension from the usual constraint $\text{hidden_dim} = \text{qhead_num} \times \text{head_dim}$ and expand only hidden-dimension; the head dimension and the numbers of QKV heads remain unchanged.

which applies copy/random/zero initialization without any rescaling, and (2) *RMS-Preserved Scaled*, which applies the rescaling derived in Sec. 3.2 to enforce activation RMS consistency at the expansion moment.

Results. Figure 1 reports expert-inner expansion, enumerating five fan-out/fan-in initialization pairs within each expert MLP, denoted as *up_proj – down_proj init*.

Across all initializations, *Naive Init*, *No Scaled* consistently yields a smaller immediate loss gap but worse late-stage convergence, whereas *RMS-Preserved Scaled* recovers steadily and converges to a lower final loss. The hidden-dimension expansion counterpart in Appendix D shows the same pattern. Overall, RMS-preserving expansion robustly improves late-stage convergence under both expert-inner and hidden-dimension growth across diverse initialization strategies.

Figure 1(f) aggregates late-stage performance across initialization pairs under *RMS-Preserved Scaled*. While broadly beneficial, the *both-sides copied* configuration still significantly underperforms the other RMS-preserved variants.

4 Breaking the Symmetry Lock

The experimental results in Sec. 3.3 reveal a counter-intuitive phenomenon: although the copy strategy strictly preserves the forward output at expansion, it consistently underperforms other RMS-preserved initializations, exhibiting both slower post-expansion recovery and a higher eventual loss. Intuitively, copy-based initialization seems ideal, as it ensures a seamless loss transition and thus the most stable starting point. We argue that the gap is instead governed by a copy-induced backward-pass symmetry: duplicated components receive identical gradients and thus evolve identically, failing to diversify into distinct features and rendering the expanded capacity functionally redundant. We formally derive this mechanism in the following analysis.

4.1 Identical Gradients Under Copy Expansion

Consider the linear layer in Eq. (2). We analyze gradient dynamics under $2\times$ width expansion with copy initialization.

Fan-Out Expansion. Specializing Eq. (8) to copy initialization gives $\tilde{\mathbf{W}} = \mathbf{W}$ and $\mathbf{W}' = [\mathbf{W}^\top, \mathbf{W}^\top]^\top$, hence $\mathbf{y}' = [\mathbf{y}, \mathbf{y}]$. If subsequent layers are also copied, back-propagation maintains symmetry: $\frac{\partial \mathcal{L}}{\partial \mathbf{y}'} = [\mathbf{g}, \mathbf{g}]$ with $\mathbf{g} = \frac{\partial \mathcal{L}}{\partial \mathbf{y}}$. The gradient w.r.t. the expanded weights is:

$$\nabla_{\mathbf{W}'} \mathcal{L} = \frac{\partial \mathcal{L}}{\partial \mathbf{y}'} \mathbf{x}^\top = \begin{bmatrix} \mathbf{g} \\ \mathbf{g} \end{bmatrix} \mathbf{x}^\top = \begin{bmatrix} \mathbf{g} \mathbf{x}^\top \\ \mathbf{g} \mathbf{x}^\top \end{bmatrix}. \quad (20)$$

We provide the analogous analysis for fan-in expansion in Appendix A.3.

Symmetry Lock. In both cases, $\nabla_{\mathbf{W}} \mathcal{L} = \nabla_{\tilde{\mathbf{W}}} \mathcal{L}$ holds, indicating identical gradients in copy expansion. With symmetrically initialized optimizer states, i.e., identical momentum for AdamW, the two blocks receive identical updates, enforcing $\mathbf{W}(t) = \tilde{\mathbf{W}}(t)$ throughout training. This creates a “symmetry lock”: despite increased parameters, the model remains in the original lower-dimensional subspace. The expanded neurons fail to learn distinct features, making width scaling inefficient unless the symmetry is explicitly broken.

Orthogonalization Fails to Break Symmetry. Advanced optimizers like Muon attempt to decorrelate updates by applying Newton-Schulz orthogonalization to the matrix-valued momentum, yet this mechanism fails to break the symmetry under copy-based expansion. Importantly, this step is typically implemented as a polynomial map of the Gram matrix: for a matrix \mathbf{X}_k , the next iterate can be written in the generic form

$$\mathbf{X}_{k+1} = \mathbf{X}_k \phi(\mathbf{X}_k^\top \mathbf{X}_k), \quad (21)$$

where $\phi(\cdot)$ is a matrix polynomial corresponds to $\phi(\mathbf{G}) = \frac{1}{2}(3\mathbf{I} - \mathbf{G})$ or higher-order variants $\phi(\mathbf{G}) = \alpha\mathbf{I} + \beta\mathbf{G} + \gamma\mathbf{G}^2$ with appropriate coefficients [15].

Consider the column-duplicated (fan-in) case where the momentum is initialized as $\mathbf{X}_0 = [\mathbf{A}_0, \mathbf{A}_0]$. Let $\mathbf{P}_k = \mathbf{A}_k^\top \mathbf{A}_k$. Then the Gram matrix remains block-constant:

$$\mathbf{X}_k^\top \mathbf{X}_k = \begin{bmatrix} \mathbf{A}_k^\top \\ \mathbf{A}_k^\top \end{bmatrix} [\mathbf{A}_k, \mathbf{A}_k] = \begin{bmatrix} \mathbf{P}_k & \mathbf{P}_k \\ \mathbf{P}_k & \mathbf{P}_k \end{bmatrix}. \quad (22)$$

Thus, applying the generic orthogonalization update yields

$$\mathbf{X}_{k+1} = [\mathbf{A}_k, \mathbf{A}_k] \phi \left(\begin{bmatrix} \mathbf{P}_k & \mathbf{P}_k \\ \mathbf{P}_k & \mathbf{P}_k \end{bmatrix} \right) := [\mathbf{A}_{k+1}, \mathbf{A}_{k+1}]. \quad (23)$$

Since $\phi(\cdot)$ preserves block-exchange symmetry, the update in Eq. (23) retains two identical column blocks. Therefore, the orthogonalization step cannot spontaneously break the symmetry lock induced by copy initialization.

4.2 Breaking Symmetry in Practice

4.2.1 Optimizer State Reset as a Necessary Intervention

Copy-based expansion yields identical gradients for the original and duplicated parameters. If the optimizer states are also initialized symmetrically—either by copying the existing states or by resetting all states to zero—the two halves receive identical updates, so the symmetry lock persists under both AdamW and Muon. To break this coupling symmetry without discarding the original model’s training signal, we enforce an *asymmetric* treatment: retaining the optimizer states for the original \mathbf{W} and resetting the states for the new parameters $\tilde{\mathbf{W}}$, the corresponding optimizer state matrix \mathbf{S}' (representing both first and second momentum for AdamW and momentum for Muon)

$$\mathbf{S}' = [\mathbf{S}, \mathbf{0}], \quad (24)$$

where \mathbf{S} is the pre-expansion state of \mathbf{W} and $\mathbf{0}$ initializes the state of $\tilde{\mathbf{W}}$.

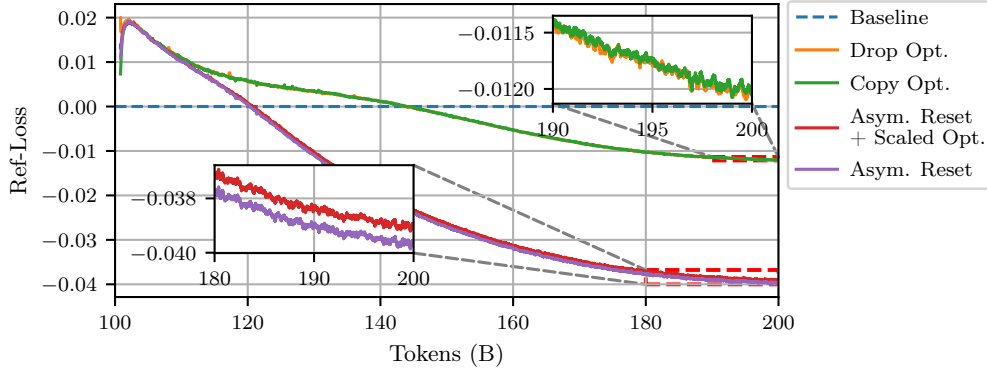


Figure 2 Optimizer-state handling under copy-based expansion. Symmetric treatments (*Drop/Copy*) exhibit a symmetry lock, yielding slower recovery and higher loss. Our asymmetric reset avoids this bottleneck, while state scaling provides no additional gain.

Experimental setup. Following Sec. 3.3, we study expert-inner expansion under copy-copy initialization and vary only the optimizer-state handling at expansion, while applying RMS-preserving parameter scaling in all variants to keep forward activation scales consistent. We compare four treatments: (i) *Drop Opt.*, globally reset all states, (ii) *Copy Opt.*, duplicate states, (iii) *Asymmetric Reset*, reset states only for new channels, and (iv) *Asymmetric Reset + Scaled Opt.*, additionally applying the parameter scaling to optimizer states. The first two are symmetric, whereas the latter two break symmetry via optimizer dynamics.

Results. Figure 2 reveals a clear gap between symmetric and asymmetric treatments. Both symmetric baselines (*Drop Opt.* and *Copy Opt.*) underperform substantially, showing slower post-expansion recovery and a higher converged loss, consistent with copy-induced backward symmetry that keeps duplicated components

tightly coupled. In contrast, *Asymmetric Reset* improves both recovery speed and final loss, indicating that resetting optimizer states only for the new channels suffices to break the symmetry lock and enable feature diversification. Notably, explicit optimizer-state scaling (*Asymmetric Reset + Scaled Opt.*) yields no additional gain, suggesting that strict state-parameter scale alignment is unnecessary and any initial mis-scaling can be quickly corrected by subsequent gradient updates.

4.2.2 Asymmetric Learning Rate Re-warmup

For training from scratch, we use a standard cosine decay learning rate scheduler with linear warmup. Let T_w denote the number of re-warmup steps, T the total number of steps, and let η_0, η_{\max} and η_{\min} be the initial, peak and final learning rates, respectively. The baseline schedule is

$$\eta(t) = f(t; T_w, T, \eta_0, \eta_{\max}, \eta_{\min}) = \begin{cases} \eta_0 + (\eta_{\max} - \eta_0) \cdot \frac{t}{T_w}, & 0 \leq t < T_w, \\ \eta_{\min} + (\eta_{\max} - \eta_{\min}) \psi\left(\frac{t - T_w}{T - T_w}\right), & T_w \leq t \leq T, \end{cases} \quad (25)$$

where

$$\psi(x) = \frac{1}{2} (1 + \cos(\pi x)). \quad (26)$$

At an expansion point t_e , we keep the original parameters on the same baseline schedule to preserve continuity, i.e., $\eta(t) = f(t; T_w, T, \eta_0, \eta_{\max}, \eta_{\min})$ for all t .

For the newly introduced parameters, we perform an asymmetric re-warmup that starts exactly from the current learning rate $\eta_e = \eta(t_e)$ and warms up for τ_w steps to a new peak learning rate proportional to η_e :

$$\hat{\eta}_{\max} = \rho \cdot \eta_e, \quad \eta_e = \eta(t_e), \quad (27)$$

where ρ is the rewarmup ratio. The learning rate for the new parameters is then defined as

$$\eta_{\text{new}}(t) = f(t - t_e; \tau_w, T - t_e, \eta_e, \hat{\eta}_{\max}, \eta_{\min}), \quad t > t_e, \quad (28)$$

where, after rewarmup, the schedule follows the same cosine-decay regime and decays to the shared minimum learning rate η_{\min} . See Appendix E for a sample curve.

4.3 Asymmetric Learning Rate Re-Warmup Further Improves Convergence Consistently.

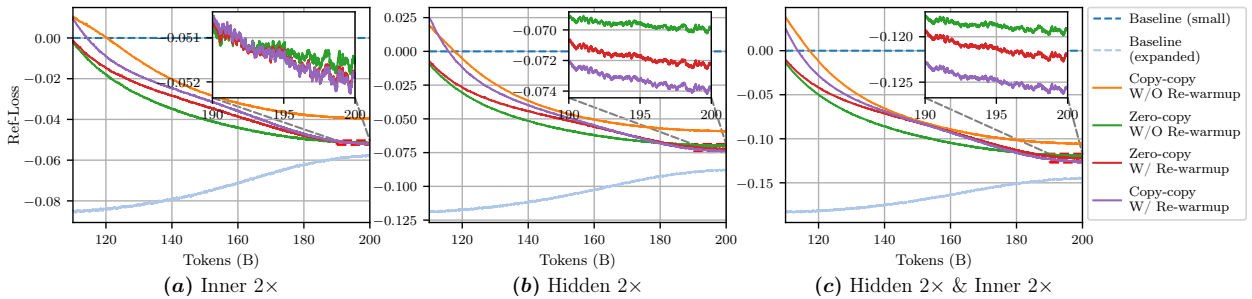


Figure 3 Asymmetric re-warmup consistently improves convergence under mid-stage width expansion.

Across Inner 2 \times , Hidden 2 \times , and joint expansion, re-warmup lowers the final loss for both RMS-preserved copy-copy and zero-copy. Copy-copy benefits most, achieving the best final loss, effectively mitigating copy-induced symmetry lock.

Experimental setup. Following Sec. 3.3, we evaluate *asymmetric learning rate re-warmup* across different width-growth axes. We consider three expansion settings: expert-inner (Inner 2 \times), hidden-dimension (Hidden

$2\times$), and joint (*Hidden* $2\times$ & *Inner* $2\times$). For all settings, we apply RMS-preserved scaling and asymmetric optimizer-state resetting, then ablate re-warmup by comparing runs with vs. without it. We report both copy-copy and zero-copy initializations (the best-performing no-re-warmup setting in our earlier analysis, see Figure 1) to assess robustness. We set the re-warmup ratio $\rho = 1.3$ and the number of re-warmup steps $\tau_w = 250$ based on Appendix F.

Results. Figure 3 shows that asymmetric learning rate re-warmup consistently improves convergence across width axes and initialization strategies. Across all width axes in three settings, enabling re-warmup yields a lower eventual loss under the same token budget. The benefit holds for both zero-copy, where new channels begin with near-zero forward contribution, and copy-copy, which strictly preserves the forward mapping. Notably, the gain is largest for copy-copy: re-warmup closes the post-expansion gap with zero-copy and reaches the lowest final loss among variants. Consistent with Sec. 4.1, beyond RMS preservation and asymmetric state reset, re-warmup adds controlled optimization asymmetry that encourages duplicated subspaces to diversify into effective capacity. We further verify SPARKLING’s generality across optimizer families (e.g., Muon) and superiority over heuristic baselines. See Appendix G and H for details. Overall, re-warmup is a robust component of SPARKLING, reliably improving convergence across width-growth axes and initialization regimes.

5 Discussions

Taken together, our SPARKLING framework comprises (i) RMS-preserved scaling, (ii) copy-based initialization, (iii) asymmetric optimizer-state reset, and (iv) asymmetric learning rate re-warmup schedule. In this section, we evaluate its overall performance.

5.1 Overall Downstream Performance

Table 1 Downstream performance under $2\times$ mid-stage width growth. Across Inner $2\times$, Hidden $2\times$, and joint expansion, SPARKLING matches or outperforms the from-scratch expanded baseline on most tasks and achieves the best average, despite a slightly higher final pre-training loss.

Model	Loss (\downarrow)	ARC-C (\uparrow)	ARC-E (\uparrow)	Arith. (\uparrow)	BoolQ (\uparrow)	CSQA (\uparrow)	HellaS. (\uparrow)	MMLU (\uparrow)	OBQA (\uparrow)	PIQA (\uparrow)	SciQ (\uparrow)	SIQA (\uparrow)	WinoG. (\uparrow)	Avg. (\uparrow)
Baseline (small)	2.3673	41.47	72.46	43.63	66.36	47.58	65.86	32.75	39.40	76.28	92.50	46.57	62.19	57.26
<i>Inner $2\times$</i>														
Baseline (expand)	2.3096	43.14	74.56	55.67	67.80	47.58	69.45	32.67	42.40	78.18	92.80	46.67	64.80	59.64
Naive Expansion	2.3276	44.15	73.86	51.10	66.45	48.24	68.04	32.71	41.80	77.09	92.90	46.98	64.25	58.96
SPARKLING	2.3153	43.48	74.39	60.50	66.82	49.06	69.21	34.05	41.60	78.35	93.30	47.80	65.04	60.30
<i>Hidden $2\times$</i>														
Baseline (expand)	2.2795	44.48	77.72	54.47	67.06	48.32	70.66	33.82	42.00	78.56	93.90	47.85	66.61	60.46
Naive Expansion	2.3082	41.81	73.86	53.77	67.37	49.06	69.40	32.00	41.00	78.18	93.00	47.44	64.17	59.25
SPARKLING	2.2933	44.48	76.14	61.90	67.16	49.80	70.06	34.49	43.40	78.45	93.40	48.26	65.98	61.13
<i>Hidden $2\times$ & Inner $2\times$</i>														
Baseline (expand)	2.2225	46.82	76.14	55.03	67.65	49.71	72.54	33.20	43.80	79.00	93.30	48.57	64.88	60.89
Naive Expansion	2.2615	45.82	74.56	58.67	67.09	51.60	71.92	33.64	41.40	79.00	93.70	47.54	65.75	60.89
SPARKLING	2.2415	46.82	77.19	66.00	69.08	50.86	72.82	35.07	43.00	79.11	94.10	48.36	68.19	62.55

Experimental setup. Following Sec. 4.3, we further evaluate downstream performance under three $2\times$ width-growth settings. For each setting, we compare (i) the *Baseline (small)* model before expansion, (ii) the *Baseline (expand)* model trained from scratch at the target width under the same token budget, (iii) *Naive Expansion* variant with copy-based initialization but without our interventions, and (iv) our *SPARKLING*, which combines RMS-preserved scaling with asymmetric optimizer-state reset and asymmetric learning rate re-warmup for newly introduced parameters. We report the final pre-training loss and downstream accuracies, including ARC-C/E [8], Arithmetic [4], BoolQ [7], CommonsenseQA [29], HellaSwag [43], MMLU [14], OpenBookQA [21], PIQA [3], SciQ [32], SocialIQA [26], Winogrande [25].

Results. Table 1 shows that mid-stage expansion still leaves a small gap in final pre-training loss relative to training the expanded model from scratch. Nevertheless, SPARKLING achieves the best downstream average among the expansion variants and matches or exceeds the from-scratch expanded baseline on most tasks. Overall, these results validate the reliability of SPARKLING: despite slightly larger final pre-training loss, our framework consistently improves downstream performance across diverse width-growth axes.

5.2 Computational Cost Analysis

Table 2 Compute-cost comparison under a fixed token budget.

Method	Total Tokens	Expand @	Act./Tot. Params	FLOPs ($\times 10^{20}$)	Wall- clock (h)	FLOPs Saved	Speed -up
Baseline	200B	-	450M/2.56B	5.40	48		
<i>Inner 2×</i>							
From Scratch SPARKLING	200B	- 100B	751M/5B	9.01 7.21	84 66	- 20%	- 1.27×
<i>Hidden 2×</i>							
From Scratch SPARKLING	200B	- 100B	900M/5.13B	10.80 8.10	96 75	- 25%	- 1.29×
<i>Hidden 2×</i> & <i>Inner 2×</i>							
From Scratch SPARKLING	200B	- 100B	1.5B/9.96B	18.00 11.70	209 140	- 35%	- 1.49×

Now that the effectiveness of our expansion framework has been validated, we finally return to the core motivation of progressive learning - reducing training costs while retaining the performance of the target-width model. We quantify the computational savings by comparing mid-stage width expansion with training the target-width model from scratch *under the same training token budget*. Following Kaplan et al. [16], we approximate pre-training compute as $C \approx 6ND$, where N is the number of *active* parameters and D is the total training tokens. Suppose that the expansion occurs at D_e tokens, where the small model with active size N_{small} is trained for D_e tokens, and the expanded model of active size N_{large} is trained for $(D - D_e)$ tokens, yielding

$$C^* \approx 6(N_{\text{small}}D_e + N_{\text{large}}(D - D_e)), \quad (29)$$

whereas training the expanded model from scratch costs $C_{\text{scratch}} \approx 6N_{\text{large}}D$. We report the relative reduction as $\text{FLOPs Saved} = 1 - C^*/C_{\text{scratch}}$, and the empirical wall-clock *Speed-up* as T_{scratch}/T^* .

Table 2 summarizes results across three $2\times$ width-growth settings. Under the same 200B-token budget, SPARKLING saves 20%–35% training FLOPs relative to training the expanded model from scratch, and achieves up to a $1.49\times$ measured wall-clock speed-up under $2\times$ width expansion. Overall, SPARKLING matches or even exceeds the performance of the from-scratch expanded model while substantially reducing training costs, making mid-stage width growth practically advantageous.

6 Conclusion and Future Work

We proposed SPARKLING, a systematic progressive learning framework via width expansion, and resolved the challenges arising during mid-stage model expansion. In contrast to conventional function-preserving perspectives, we emphasize signal preservation by maintaining the RMS-scale of activations during expansion. To break the symmetry induced by copy-based initialization, we apply asymmetric optimizer resetting together with learning rate re-warmup. Across multiple width axes and optimizer families, extensive experiments validate both the effectiveness and the efficiency of our framework.

While our results are promising, several avenues remain for future exploration. First, a unified principle of simultaneous width and depth expansion has yet to be established. Moreover, we aim to investigate whether our RMS preservation strategy could satisfy the μP condition [38], where the transferability of optimal hyperparameters is naturally ensured after expansion. We view these as critical future work toward developing a more comprehensive, “tuning-free” framework for progressive learning.

References

- [1] Thomas Bachlechner, Bodhisattwa Prasad Majumder, Huanru Henry Mao, Garrison W. Cottrell, and Julian McAuley. Rezero is all you need: Fast convergence at large depth, 2020. URL <https://arxiv.org/abs/2003.04887>.
- [2] George Bird. The affine divergence: Aligning activation updates beyond normalisation, 2025. URL <https://arxiv.org/abs/2512.22247>.
- [3] Yonatan Bisk, Rowan Zellers, Ronan Le bras, Jianfeng Gao, and Yejin Choi. Piqa: Reasoning about physical commonsense in natural language. *Proceedings of the AAAI Conference on Artificial Intelligence*, 34(05):7432–7439, Apr. 2020. doi: 10.1609/aaai.v34i05.6239. URL <https://ojs.aaai.org/index.php/AAAI/article/view/6239>.
- [4] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020.
- [5] Cheng Chen, Yichun Yin, Lifeng Shang, Xin Jiang, Yujia Qin, Fengyu Wang, Zhi Wang, Xiao Chen, Zhiyuan Liu, and Qun Liu. bert2BERT: Towards reusable pretrained language models. In Smaranda Muresan, Preslav Nakov, and Aline Villavicencio, editors, *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 2134–2148, Dublin, Ireland, May 2022. Association for Computational Linguistics. doi: 10.18653/v1/2022.acl-long.151. URL <https://aclanthology.org/2022.acl-long.151/>.
- [6] Tianqi Chen, Ian Goodfellow, and Jonathon Shlens. Net2net: Accelerating learning via knowledge transfer, 2016. URL <https://arxiv.org/abs/1511.05641>.
- [7] Christopher Clark, Kenton Lee, Ming-Wei Chang, Tom Kwiatkowski, Michael Collins, and Kristina Toutanova. BoolQ: Exploring the surprising difficulty of natural yes/no questions. In Jill Burstein, Christy Doran, and Tamar Solorio, editors, *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 2924–2936, Minneapolis, Minnesota, June 2019. Association for Computational Linguistics. doi: 10.18653/v1/N19-1300. URL <https://aclanthology.org/N19-1300/>.
- [8] Peter Clark, Isaac Cowhey, Oren Etzioni, Tushar Khot, Ashish Sabharwal, Carissa Schoenick, and Oyvind Tafjord. Think you have solved question answering? try arc, the ai2 reasoning challenge. *arXiv preprint arXiv:1803.05457*, 2018.
- [9] DeepSeek-AI, Aixin Liu, Bei Feng, Bing Xue, Bingxuan 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, Fuli Luo, Guangbo Hao, Guanting Chen, Guowei Li, H. Zhang, Han Bao, Hanwei Xu, Haocheng Wang, Haowei Zhang, Honghui Ding, Huajian Xin, Huazuo Gao, Hui Li, Hui Qu, J. L. Cai, Jian Liang, Jianzhong Guo, Jiaqi Ni, Jiashi Li, Jiawei Wang, Jin Chen, Jingchang Chen, Jingyang Yuan, Junjie Qiu, Junlong Li, Junxiao Song, Kai Dong, Kai Hu, Kaige Gao, Kang Guan, Kexin Huang, Kuai Yu, Lean Wang, Lecong Zhang, Lei Xu, Leyi Xia, Liang Zhao, Litong Wang, Liyue Zhang, Meng Li, Miaojuan Wang, Mingchuan Zhang, Minghua Zhang, Minghui Tang, Mingming Li, Ning Tian, Panpan Huang, Peiyi Wang, Peng Zhang, Qiancheng Wang, Qihao Zhu, Qinyu Chen, Qiushi Du, R. J. Chen, R. L. Jin, Ruiqi Ge, Ruisong Zhang, Ruizhe Pan, Runji Wang, Runxin Xu, Ruoyu Zhang, Ruyi Chen, S. S. Li, Shanghao Lu, Shangyan Zhou, Shanhuang Chen, Shaoqing Wu, Shengfeng Ye, Shengfeng Ye, Shirong Ma, Shiyu Wang, Shuang Zhou, Shuiping Yu, Shunfeng Zhou, Shuting Pan, T. Wang, Tao Yun, Tian Pei, Tianyu Sun, W. L. Xiao, Wangding Zeng, Wanbiao Zhao, Wei An, Wen Liu, Wenfeng Liang, Wenjun Gao, Wenqin Yu, Wentao Zhang, X. Q. Li, Xiangyue Jin, Xianzu Wang, Xiao Bi, Xiaodong Liu, Xiaohan Wang, Xiaojin Shen, Xiaokang Chen, Xiaokang Zhang, Xiaosha Chen, Xiaotao Nie, Xiaowen Sun, Xiaoxiang Wang, Xin Cheng, Xin Liu, Xin Xie, Xingchao Liu, Xingkai Yu, Xinnan Song, Xinxia Shan, Xinyi Zhou, Xinyu Yang, Xinyuan Li, Xuecheng Su, Xuheng Lin, Y. K. Li, Y. Q. Wang, Y. X. Wei, Y. X. Zhu, Yang Zhang, Yanhong Xu, Yanhong Xu, Yanping Huang, Yao Li, Yao Zhao, Yaofeng Sun, Yaohui Li, Yaohui Wang, Yi Yu, Yi Zheng, Yichao Zhang, Yifan Shi, Yiliang Xiong, Ying He, Ying Tang, Yishi Piao, Yisong Wang, Yixuan Tan, Yiyang Ma, Yiyuan Liu, Yongqiang Guo, Yu Wu, Yuan Ou, Yuchen Zhu, Yudian Wang, Yue Gong, Yuheng Zou, Yujia He, Yukun Zha, Yunfan Xiong, Yunxian Ma, Yuting Yan, Yuxiang Luo, Yuxiang You, Yuxuan Liu, Yuyang Zhou, Z. F. Wu, Z. Z. Ren, Zehui Ren, Zhangli Sha, Zhe Fu, Zhean Xu, Zhen Huang, Zhen Zhang, Zhenda Xie, Zhengyan Zhang, Zhewen Hao, Zhibin Gou, Zhicheng Ma, Zhigang Yan, Zhihong Shao, Zhipeng Xu, Zhiyu Wu, Zhongyu Zhang, Zhuoshu Li, Zihui Gu, Zijia Zhu, Zijun Liu, Zilin Li, Ziwei Xie, Ziyang Song, Ziyi Gao, and Zizheng Pan. Deepseek-v3 technical report, 2025. URL <https://arxiv.org/abs/2412.19437>.

- [10] Wenyu Du, Tongxu Luo, Zihan Qiu, Zeyu Huang, Yikang Shen, Reynold Cheng, Yike Guo, and Jie Fu. Stacking your transformers: A closer look at model growth for efficient llm pre-training. *Advances in Neural Information Processing Systems*, 37:10491–10540, 2024.
- [11] Utku Evci, Bart van Merriënboer, Thomas Unterthiner, Max Vladymyrov, and Fabian Pedregosa. Gradmax: Growing neural networks using gradient information, 2022. URL <https://arxiv.org/abs/2201.05125>.
- [12] Linyuan Gong, Di He, Zhuohan Li, Tao Qin, Liwei Wang, and Tieyan Liu. Efficient training of bert by progressively stacking. In *International conference on machine learning*, pages 2337–2346. PMLR, 2019.
- [13] Xue Han, Yitong Wang, Junlan Feng, Qian Hu, Chao Deng, et al. Loire: Lifelong learning on incremental data via pre-trained language model growth efficiently. In *The Thirteenth International Conference on Learning Representations*, 2025.
- [14] Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. Measuring massive multitask language understanding. In *International Conference on Learning Representations*, 2021. URL <https://openreview.net/forum?id=d7KBjmI3GmQ>.
- [15] Keller Jordan, Yuchen Jin, Vlado Boza, Jiacheng You, Franz Cesista, Laker Newhouse, and Jeremy Bernstein. Muon: An optimizer for hidden layers in neural networks, 2024. URL <https://kellerjordan.github.io/posts/muon/>.
- [16] Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B. Brown, Benjamin Chess, Rewon Child, Scott Gray, Alec Radford, Jeffrey Wu, and Dario Amodei. Scaling laws for neural language models, 2020. URL <https://arxiv.org/abs/2001.08361>.
- [17] Sanghoon Kim, Dahyun Kim, Chanjun Park, Wonsung Lee, Wonho Song, Yunsu Kim, Hyeonwoo Kim, Yungi Kim, Hyeonju Lee, Jihoo Kim, et al. Solar 10.7 b: Scaling large language models with simple yet effective depth up-scaling. In *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 6: Industry Track)*, pages 23–35, 2024.
- [18] Houyi Li, Wenzhen Zheng, Qiufeng Wang, Hanshan Zhang, Zili Wang, Shijie Xuyang, Yuantao Fan, Zhenyu Ding, Haoying Wang, Ning Ding, Shuigeng Zhou, Xiangyu Zhang, and Daxin Jiang. Predictable scale: Part i, step law – optimal hyperparameter scaling law in large language model pretraining, 2025. URL <https://arxiv.org/abs/2503.04715>.
- [19] Jingyuan Liu, Jianlin Su, Xingcheng Yao, Zhejun Jiang, Guokun Lai, Yulun Du, Yidao Qin, Weixin Xu, Enzhe Lu, Junjie Yan, Yanru Chen, Huabin Zheng, Yibo Liu, Shaowei Liu, Bohong Yin, Weiran He, Han Zhu, Yuzhi Wang, Jianzhou Wang, Mengnan Dong, Zheng Zhang, Yongsheng Kang, Hao Zhang, Xinran Xu, Yutao Zhang, Yuxin Wu, Xinyu Zhou, and Zhilin Yang. Muon is scalable for llm training, 2025. URL <https://arxiv.org/abs/2502.16982>.
- [20] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. In *International Conference on Learning Representations*, 2019. URL <https://openreview.net/forum?id=Bkg6RiCqY7>.
- [21] Todor Mihaylov, Peter Clark, Tushar Khot, and Ashish Sabharwal. Can a suit of armor conduct electricity? a new dataset for open book question answering. In Ellen Riloff, David Chiang, Julia Hockenmaier, and Jun’ichi Tsujii, editors, *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 2381–2391, Brussels, Belgium, October–November 2018. Association for Computational Linguistics. doi: 10.18653/v1/D18-1260. URL <https://aclanthology.org/D18-1260/>.
- [22] Niklas Muennighoff, Luca Soldaini, Dirk Groeneveld, Kyle Lo, Jacob Morrison, Sewon Min, Weijia Shi, Evan Pete Walsh, Oyvind Tafjord, Nathan Lambert, Yuling Gu, Shane Arora, Akshita Bhagia, Dustin Schwenk, David Wadden, Alexander Wettig, Binyuan Hui, Tim Dettmers, Douwe Kiela, Ali Farhadi, Noah A. Smith, Pang Wei Koh, Amanpreet Singh, and Hannaneh Hajishirzi. OLMoe: Open mixture-of-experts language models. In *The Thirteenth International Conference on Learning Representations*, 2025. URL <https://openreview.net/forum?id=xXTkbTBmqj>.
- [23] Niklas Muennighoff, Luca Soldaini, Dirk Groeneveld, Kyle Lo, Jacob Morrison, Sewon Min, Weijia Shi, Pete Walsh, Oyvind Tafjord, Nathan Lambert, Yuling Gu, Shane Arora, Akshita Bhagia, Dustin Schwenk, David Wadden, Alexander Wettig, Binyuan Hui, Tim Dettmers, Douwe Kiela, Ali Farhadi, Noah A. Smith, Pang Wei Koh, Amanpreet Singh, and Hannaneh Hajishirzi. Olmoe: Open mixture-of-experts language models, 2025. URL <https://arxiv.org/abs/2409.02060>.
- [24] Yu Pan, Ye Yuan, Yichun Yin, Zenglin Xu, Lifeng Shang, Xin Jiang, and Qun Liu. Reusing pretrained models by multi-linear operators for efficient training. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023. URL <https://openreview.net/forum?id=RgNXKIrwYU>.

- [25] Keisuke Sakaguchi, Ronan Le Bras, Chandra Bhagavatula, and Yejin Choi. Winogrande: An adversarial winograd schema challenge at scale. *Proceedings of the AAAI Conference on Artificial Intelligence*, 34(05):8732–8740, Apr. 2020. doi: 10.1609/aaai.v34i05.6399. URL <https://ojs.aaai.org/index.php/AAAI/article/view/6399>.
- [26] Maarten Sap, Hannah Rashkin, Derek Chen, Ronan Le Bras, and Yejin Choi. Social IQa: Commonsense reasoning about social interactions. In Kentaro Inui, Jing Jiang, Vincent Ng, and Xiaojun Wan, editors, *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 4463–4473, Hong Kong, China, November 2019. Association for Computational Linguistics. doi: 10.18653/v1/D19-1454. URL <https://aclanthology.org/D19-1454/>.
- [27] Andrew M. Saxe, James L. McClelland, and Surya Ganguli. Exact solutions to the nonlinear dynamics of learning in deep linear neural networks, 2014. URL <https://arxiv.org/abs/1312.6120>.
- [28] Sheng Shen, Pete Walsh, Kurt Keutzer, Jesse Dodge, Matthew Peters, and Iz Beltagy. Staged training for transformer language models, 2022. URL <https://arxiv.org/abs/2203.06211>.
- [29] Alon Talmor, Jonathan Herzig, Nicholas Lourie, and Jonathan Berant. CommonsenseQA: A question answering challenge targeting commonsense knowledge. In Jill Burstein, Christy Doran, and Thamar Solorio, editors, *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4149–4158, Minneapolis, Minnesota, June 2019. Association for Computational Linguistics. doi: 10.18653/v1/N19-1421. URL <https://aclanthology.org/N19-1421/>.
- [30] Peihao Wang, Rameswar Panda, Lucas Torroba Hennigen, Philip Greengard, Leonid Karlinsky, Rogerio Feris, David Daniel Cox, Zhangyang Wang, and Yoon Kim. Learning to grow pretrained models for efficient transformer training. In *The Eleventh International Conference on Learning Representations*, 2023. URL <https://openreview.net/forum?id=cDYRS5iZ16f>.
- [31] Yite Wang, Jiahao Su, Hanlin Lu, Cong Xie, Tianyi Liu, Jianbo Yuan, Haibin Lin, Ruoyu Sun, and Hongxia Yang. LEMON: Lossless model expansion. In *The Twelfth International Conference on Learning Representations*, 2024. URL <https://openreview.net/forum?id=3Vw7DQqq7U>.
- [32] Johannes Welbl, Nelson F. Liu, and Matt Gardner. Crowdsourcing multiple choice science questions. In Leon Derczynski, Wei Xu, Alan Ritter, and Tim Baldwin, editors, *Proceedings of the 3rd Workshop on Noisy User-generated Text*, pages 94–106, Copenhagen, Denmark, September 2017. Association for Computational Linguistics. doi: 10.18653/v1/W17-4413. URL <https://aclanthology.org/W17-4413/>.
- [33] Chengyue Wu, Yukang Gan, Yixiao Ge, Zeyu Lu, Jiahao Wang, Ye Feng, Ying Shan, and Ping Luo. LLaMA pro: Progressive LLaMA with block expansion. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 6518–6537, 2024.
- [34] Lemeng Wu, Dilin Wang, and Qiang Liu. Splitting steepest descent for growing neural architectures. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d’Alché-Buc, E. Fox, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc., 2019. URL https://proceedings.neurips.cc/paper_files/paper/2019/file/3a01fc0853eb94fde4d1cc6fb842a-Paper.pdf.
- [35] Lemeng Wu, Bo Liu, Peter Stone, and Qiang Liu. Firefly neural architecture descent: a general approach for growing neural networks. In H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin, editors, *Advances in Neural Information Processing Systems*, volume 33, pages 22373–22383. Curran Associates, Inc., 2020. URL https://proceedings.neurips.cc/paper_files/paper/2020/file/fdbe012e2e11314b96402b32c0df26b7-Paper.pdf.
- [36] An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Gao, Chengen Huang, Chenxu Lv, et al. Qwen3 technical report. *arXiv preprint arXiv:2505.09388*, 2025.
- [37] Cheng Yang, Shengnan Wang, Chao Yang, Yuechuan Li, Ru He, and Jingqiao Zhang. Progressively stacking 2.0: A multi-stage layerwise training method for bert training speedup. *arXiv preprint arXiv:2011.13635*, 2020.
- [38] Greg Yang, Edward J. Hu, Igor Babuschkin, Szymon Sidor, Xiaodong Liu, David Farhi, Nick Ryder, Jakub Pachocki, Weizhu Chen, and Jianfeng Gao. Tensor programs v: Tuning large neural networks via zero-shot hyperparameter transfer, 2022. URL <https://arxiv.org/abs/2203.03466>.
- [39] Yifei Yang, Zouying Cao, Xinbei Ma, Yao Yao, Zhi Chen, Libo Qin, and Hai Zhao. Lesa: Learnable llm layer scaling-up. In *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 22463–22476, 2025.

- [40] Kazuki Yano, Sho Takase, Sosuke Kobayashi, Shun Kiyono, and Jun Suzuki. Efficient construction of model family through progressive training using model expansion. *CoRR*, abs/2504.00623, April 2025. URL <https://doi.org/10.48550/arXiv.2504.00623>.
- [41] Yiqun Yao, Zheng Zhang, Jing Li, and Yequan Wang. Masked structural growth for 2x faster language model pre-training. In *The Twelfth International Conference on Learning Representations*, 2024. URL <https://openreview.net/forum?id=rL7xsG1aRn>.
- [42] Xin Yuan, Pedro Henrique Pamplona Savarese, and Michael Maire. Accelerated training via incrementally growing neural networks using variance transfer and learning rate adaptation. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023. URL <https://openreview.net/forum?id=H1a7bVVnPK>.
- [43] Rowan Zellers, Ari Holtzman, Yonatan Bisk, Ali Farhadi, and Yejin Choi. HellaSwag: Can a machine really finish your sentence? In Anna Korhonen, David Traum, and Lluís Màrquez, editors, *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 4791–4800, Florence, Italy, July 2019. Association for Computational Linguistics. doi: 10.18653/v1/P19-1472. URL <https://aclanthology.org/P19-1472/>.
- [44] Biao Zhang and Rico Sennrich. Root mean square layer normalization. *Advances in neural information processing systems*, 32, 2019.
- [45] Bo-Wen Zhang, Liangdong Wang, Ye Yuan, Jijie Li, Shuhao Gu, Mengdi Zhao, Xinya Wu, Guang Liu, Chengwei Wu, Hanyu Zhao, Li Du, Yiming Ju, Quanyue Ma, Yulong Ao, Yingli Zhao, Songhe Zhu, Zhou Cao, Dong Liang, Yonghua Lin, Ming Zhang, Shunfei Wang, Yanxin Zhou, Min Ye, Xuekai Chen, Xinyang Yu, Xiangjun Huang, and Jian Yang. Aquilamoe: Efficient training for moe models with scale-up and scale-out strategies, 2024. URL <https://arxiv.org/abs/2408.06567>.

Appendix

A Derivations

A.1 Eq. (7): The Per-Coordinate Variance Under Fan-In Aggregation

We provide the intermediate steps used in Sec. 3.2.1 to derive RMS preservation for a fan-in variance constraint.

Consider the linear layer $\mathbf{y} = \mathbf{W}\mathbf{x}$ with $\mathbf{W} \in \mathbb{R}^{d_{\text{out}} \times d_{\text{in}}}$ and $\mathbf{x} \in \mathbb{R}^{d_{\text{in}}}$. For a fixed output coordinate $i \in \{1, \dots, d_{\text{out}}\}$, we write

$$y_i = \sum_{j=1}^{d_{\text{in}}} w_{ij} x_j. \quad (30)$$

Assume that, across the fan-in dimensions, the pairs $\{(w_{ij}, x_j)\}_{j=1}^{d_{\text{in}}}$ are mutually independent, and that \mathbf{W} and \mathbf{x} are independent of each other. Under these conditions, the variance of y_i decomposes additively:

$$\text{Var}(y_i) = \text{Var}\left(\sum_{j=1}^{d_{\text{in}}} w_{ij} x_j\right) = \sum_{j=1}^{d_{\text{in}}} \text{Var}(w_{ij} x_j). \quad (31)$$

If, in addition, the fan-in terms are homoscedastic and centered in the sense that $\mathbb{E}[w_{ij}] = \mathbb{E}[x_j] = 0$ and $\text{Var}(w_{ij}) = \sigma_w^2$, $\text{Var}(x_j) = \sigma_x^2$ for all j , then each product term shares the same variance $\text{Var}(w_{ij} x_j) = \sigma_w^2 \sigma_x^2$. Substituting this into Eq. (31) yields

$$\text{Var}(y_i) = \sum_{j=1}^{d_{\text{in}}} \sigma_w^2 \sigma_x^2 = d_{\text{in}} \sigma_w^2 \sigma_x^2, \quad (32)$$

which is the expression in Eq. (7) of the main text. Combined with Eq. (6), this shows that (when s_{in} is fixed) preserving the output RMS scale is equivalent to preserving $\text{Var}(y_i)$, and under the above assumptions this reduces to keeping $d_{\text{in}} \sigma_w^2 \sigma_x^2$ invariant across the expansion.

A.2 Eq. (16)–(17): RMS-Preserving Rescaling under Fan-In Expansion with Both-Sides Copied

A qualitatively different regime arises when *both* sides of the newly introduced fan-in coordinates are created by copying existing dimensions, i.e., the new columns of \mathbf{W}' and the new coordinates of \mathbf{x}' are both duplicated from the same subset of original fan-in dimensions, respectively. In this case, the independence across fan-in dimensions is violated: the copied pairs (w'_{ij}, x'_j) are no longer independent replicas but perfectly correlated duplicates of some original terms. As a result, the variance no longer decomposes as a simple sum of d'_{in} independent contributions, and the duplicated terms contribute quadratically through covariance.

Given Eq. (15) with c denoting the copy ratio, we first consider the setting $0 < c \leq 1$ where each copied dimension is duplicated exactly once. Let \mathcal{R} be the set of copied indices with $|\mathcal{R}| = c d_{\text{in}}$, and \mathcal{S} be the remaining indices with $|\mathcal{S}| = (1 - c) d_{\text{in}}$. Under one-to-one copying on both \mathbf{W} and \mathbf{x} , each duplicated dimension contributes twice with identical value, yielding

$$y'_i = \sum_{j=1}^{d'_{\text{in}}} w'_{ij} x'_j = 2 \sum_{r \in \mathcal{R}} w'_{ir} x_r + \sum_{s \in \mathcal{S}} w'_{is} x_s. \quad (33)$$

Under the same independent assumptions as above on the *original* terms, the variance of y'_i becomes

$$\begin{aligned}
\text{Var}(y'_i) &= \text{Var}\left(2 \sum_{r \in \mathcal{R}} w'_{ir} x_r + \sum_{s \in \mathcal{S}} w'_{is} x_s\right) \\
&= 4 \sum_{r \in \mathcal{R}} \text{Var}(w'_{ir} x_r) + \sum_{s \in \mathcal{S}} \text{Var}(w'_{is} x_s) \\
&= 4|\mathcal{R}| \sigma_{w'}^2 \sigma_x^2 + |\mathcal{S}| \sigma_{w'}^2 \sigma_x^2 \\
&= (4c d_{\text{in}} + (1-c) d_{\text{in}}) \sigma_{w'}^2 \sigma_x^2 \\
&= d_{\text{in}}(1+3c) \sigma_{w'}^2 \sigma_x^2.
\end{aligned} \tag{34}$$

When the input scale is kept unchanged, preserving the original variance requires rescaling the weights in the expanded layer. Let $\sigma_{w'}^2$ denote the post-rescaling weight variance; enforcing $\text{Var}(y'_i) = \text{Var}(y_i)$ gives

$$d_{\text{in}}(1+3c) \sigma_{w'}^2 \sigma_x^2 = d_{\text{in}} \sigma_w^2 \sigma_x^2 \implies \sigma_{w'}^2 = \frac{1}{1+3c} \sigma_w^2, \tag{35}$$

or equivalently,

$$w'_{ij} = \frac{1}{\sqrt{1+3c}} w_{ij}, \quad 0 < c \leq 1, \quad \forall i = 1, \dots, d_{\text{out}}, \quad j = 1, \dots, d'_{\text{in}}. \tag{36}$$

For $c > 1$ where some dimensions might be copied multiple times, the variance of y'_i becomes

$$\text{Var}(y'_i) = (1+c)^2 d_{\text{in}} \sigma_{w'}^2 \sigma_x^2. \tag{37}$$

When the input scale is kept unchanged, enforcing $\text{Var}(y'_i) = \text{Var}(y_i)$ gives

$$(1+c)^2 d_{\text{in}} \sigma_{w'}^2 \sigma_x^2 = d_{\text{in}} \sigma_w^2 \sigma_x^2 \implies \sigma_{w'}^2 = \frac{1}{(1+c)^2} \sigma_w^2, \tag{38}$$

or equivalently,

$$w'_{ij} = \frac{1}{1+c} w_{ij}, \quad c > 1, \quad \forall i = 1, \dots, d_{\text{out}}, \quad j = 1, \dots, d'_{\text{in}}. \tag{39}$$

Combining Eqs. (36) and (39) yields the final rescaling rule in Eq. (16). Substituting $c = d'_{\text{in}}/d_{\text{in}} - 1$ into these equations gives the equivalent form in Eq. (17).

A.3 Identical Gradients Under Copy Expansion for Fan-In Expansion

For the input expansion defined in Eq. (10), copy initialization sets $\tilde{\mathbf{W}} = \mathbf{W}$ such that $\mathbf{W}' = \alpha[\mathbf{W}, \mathbf{W}]$, with the input duplicated as $\mathbf{x}' = [\mathbf{x}, \mathbf{x}]$. The forward pass yields $\mathbf{y}' = \alpha(\mathbf{W}\mathbf{x} + \mathbf{W}\mathbf{x})$. During backward propagation:

$$\nabla_{\mathbf{W}'} \mathcal{L} = \frac{\partial \mathcal{L}}{\partial \mathbf{y}'} (\mathbf{x}')^\top = \mathbf{g} [\mathbf{x}^\top, \mathbf{x}^\top] = [\mathbf{g}\mathbf{x}^\top, \mathbf{g}\mathbf{x}^\top]. \tag{40}$$

showing that the two copied blocks receive identical gradients and that the uniform scalar α does not affect this symmetry argument.

B Detailed Experimental Setup

B.1 Baseline Model Configuration

We list the detailed hyperparameters and architectural configuration of pre-expansion baseline model before expansion in Table 3. In addition to these settings, we adopt a pre-norm design by inserting RMSNorm before both the attention and MLP sublayers. Moreover, within the attention block, we apply per-head q/k normalization along the head-dimension for each query and key head, i.e., normalizing the projected q and k vectors within each head over the d_{head} dimension.

Notably, since we tie the word embedding and the output projection, hidden-size expansion makes the shared matrix act as fan-out on the embedding side but fan-in on the output side. As our RMS-preserved scaling requires different factors for these two roles, we compensate the fan-in factor by multiplying the corresponding coefficient after the final output projection for this special case.

Table 3 Baseline model Configuration.

Configuration	Value
Number of Hidden Layers (L)	24
Hidden Size (d_{model})	1024
Expert Intermediate Size (d_{ffn})	512
Number of Attention Heads (n_{heads})	16
Number of Key/Value Heads (n_{kv})	4
Head Dimension (d_{head})	96
MoE Number of Experts (E)	64
MoE Top- k (k)	8
Embedding Size ($ \mathcal{V} $)	50304
Tie Word Embeddings	True
Activation Type	SwiGLU
Norm Type	RMSNorm
	Pre-norm
Positional Embedding	RoPE
Use Bias	False

B.2 Training Hyperparameters

We summarize the training hyperparameters in Table 4. Unless otherwise specified (e.g., the Muon experiments in Sec. G), we optimize with AdamW using $(\beta_1, \beta_2) = (0.9, 0.95)$, $\epsilon = 10^{-8}$, and weight decay 0.1, where we apply weight decay to all parameters including norms and embeddings. We use a cosine learning rate schedule with a linear warmup over 3% of total steps, and decay to a minimum learning rate set by a final ratio of 0.01 relative to the peak learning rate. All experiments are run on a cluster of $64 \times$ NVIDIA A100 GPUs with 80 GB memory each, using a global batch size of 768 with per-device microbatch size 3.

Table 4 Training hyperparameter configuration.

Configuration	Value
Total Training Tokens	200 B
Optimizer	AdamW
Peak Learning Rate	1.9556×10^{-3}
AdamW Betas (β_1, β_2)	(0.9, 0.95)
AdamW Epsilon ϵ	1.0×10^{-8}
Weight Decay	0.1
Decay Norm & Bias	True
Decay Embeddings	True
LR Scheduler	Cosine w/ Warmup
Warmup Steps	3% of total steps
Final LR Ratio	0.01
Max Sequence Length	4096
Global Batch Size	768
Device Microbatch Size	3
Number of GPUs	64

For the all models trained from scratch, we set the peak learning rate to the step-law optimum reported in

[18] and apply a batch-size scaling to obtain the corresponding value under our training setup in Table 5. In contrast, when expanding from a smaller model, we empirically find it more effective to keep the pre-expansion peak learning rate for training the enlarged model with the same peak LR.

Table 5 Step-law optimal learning rates [18] across model sizes and their batch-size-scaled learning rates.

Model	Total Tokens	Act. Params	Tot. Params	Step-law LR	Scaled LR
Baseline	200B	450M	2.56B	1.033e-3	1.449e-3
Inner 2×	200B	751M	5B	6.410e-4	8.988e-4
Hidden 2×	200B	900M	5.13B	5.760e-4	8.630e-4
Hidden 2× & Inner 2×	200B	1.5B	9.96B	3.920e-4	5.497e-4

C RMS Scale Under Zero Initialization

In Figure 4, we analyze a subtle but practically important corner case for RMS-preserving expansion: one-sided zero initialization. We consider two representative regimes, *random-zero* and *zero-copy*, and observe the RMS scale of the output and input activations of the whole MLP according to Eq. 3.

We empirically find that, when applying RMS-preserved scaling under zero initialization, the zero-initialized side should be treated as *random* rather than as a special perfectly *loss preserving* case. Intuitively, a zero-initialized block becomes a gradient-driven random distribution after the very first update, so its effective statistics quickly resemble those of a randomly initialized block. While omitting RMS-preserved scaling can be strictly loss-preserving at the expansion moment and therefore naturally satisfies RMS-scale preservation at $t = t_e$, we observe that the RMS-preserved scaling variant that treats the zero side as random yields an activation RMS ratio that remains closer to the original baseline scale as post-expansion training proceeds. In contrast, the RMS scale under naive unscaled zero initialization drifts and does not exhibit a recovery trend toward the baseline. This behavior is consistent with the fact that zero initialization necessarily disrupts the pre-expansion parameter distribution and thus requires a nontrivial number of steps to re-enter a compatible statistical regime. Accordingly, our emphasis is on the RMS scale shape after the model has taken a small number of post-expansion updates, rather than on the degenerate preservation at the boundary itself.

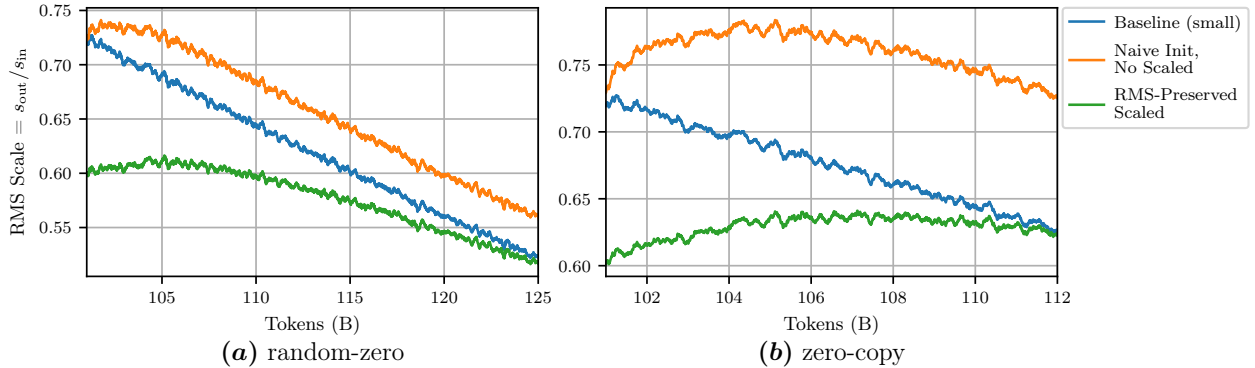


Figure 4 Under zero initialization, RMS-preserved scaling enables the post-expansion activation RMS scale to quickly recover toward the original baseline, indicating that zero initialization should be treated as random under RMS-preserving expansion.

D Hidden-Dimension Expansion: RMS-Preserved Scaling

Following the experimental setting in Sec. 3.3, we provide the hidden-dimension counterpart of our RMS-scale analysis by doubling the model hidden size from 1024 to 2048 at 100 B tokens and continuing training to 200 B tokens under the same training recipe. Figure 5 shows the same qualitative conclusion as expert-inner growth:

while naive unscaled initialization can exhibit a smaller instantaneous loss discontinuity at the expansion moment, enforcing RMS-scale consistency via our RMS-preserved rescaling yields consistently better late-stage recovery and a lower converged loss across initialization regimes.

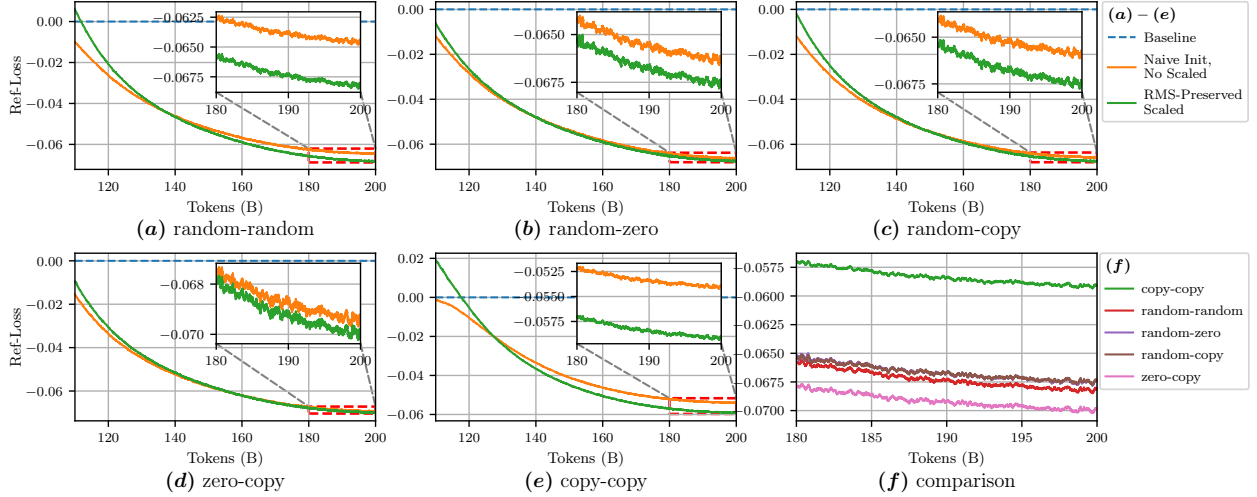


Figure 5 Hidden-dimension expansion mirrors expert-inner growth. We repeat the RMS-preserving scaling comparison under hidden-dimension $2\times$ expansion ($1024 \rightarrow 2048$ at 100 B tokens). Across initialization pairs, RMS-preserved rescaling consistently improves late-stage convergence relative to naive unscaled expansion, exhibiting the same pattern observed for expert-inner growth in Sec. 3.3.

E A Sample Asymmetric Re-warmup Learning Rate Curve

To make the asymmetric re-warmup schedule in Eq. (28) more tangible, we plot a representative learning rate trajectory in Figure 6. Typically, at the expansion point, the original parameters retrain the unchanged baseline cosine schedule for continuity, while the newly introduced parameters are re-warmed up from the current learning rate to a slightly higher peak for a short window following with decay.

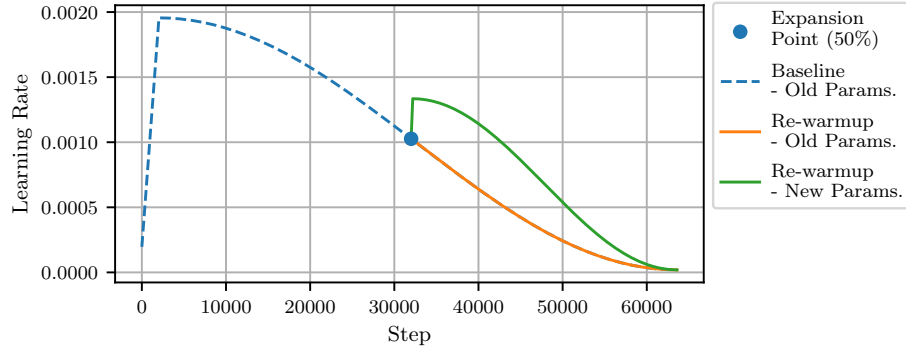


Figure 6 A sample asymmetric learning rate re-warmup curve. At the expansion step t_e , the learning rate of the original parameters remains on the baseline cosine schedule, whereas the newly introduced parameters are re-warmed from the instantaneous rate $\eta_e = \eta(t_e)$ to a higher peak $\hat{\eta}_{\max} = \rho \eta_e$ over τ_w steps, and then decay with the same cosine tail toward η_{\min} , as specified in Eq. (28).

F Hyperparameter Search for Asymmetric Re-warmup

We study how the asymmetric re-warmup schedule in Eq. (28) depends on two hyperparameters: the re-warmup ratio ρ and the number of re-warmup steps τ_w . We conduct this search under the expert-inner $2\times$ expansion

setting.

Figure 7 summarizes the final loss obtained by sweeping ρ and τ_w . The results exhibit a broad, stable region in which re-warmup is most effective: $\rho \approx 1.25$ – 1.3 and $\tau_w \approx 0$ – 250 steps achieve the lowest final loss, indicating that newly introduced parameters benefit from a modest, short-lived learning rate boost rather than a prolonged or overly strong re-warmup. Empirically, we find that this setting is also suitable for hidden-dimension expansion, and we therefore set $\rho = 1.3$ and $\tau_w = 250$ as the default hyperparameters for all experiments involving re-warmup.

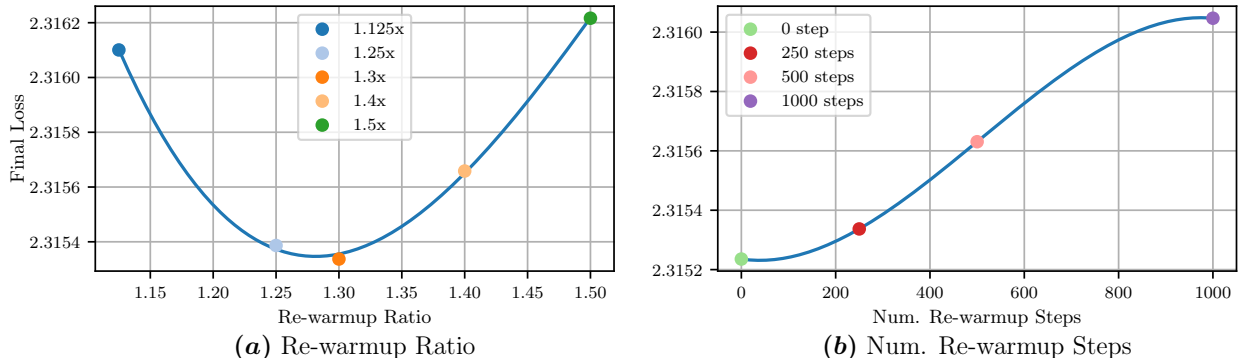


Figure 7 Hyperparameter search for asymmetric re-warmup. Under the expert-inner $2\times$ expansion setting, $\rho \approx 1.25$ – 1.3 , $\tau_w \approx 0$ – 250 , yields the lowest final loss and we adopt $\rho = 1.3$, $\tau_w = 250$ as the default re-warmup configuration in all experiments involving re-warmup.

G Effectiveness Under Muon

Experimental setup. To verify that our framework is not tied to element-wise optimizers like AdamW, we repeat the expert-inner expansion experiment following Sec. 3.3 & 4.3 using Muon as the optimizer while keeping the other recipe unchanged. We evaluate two representative components of our method. First, we isolate *RMS-preserved scaling* by comparing against the naive unscaled initialization under the same initialization regime. Second, we evaluate *asymmetric learning rate re-warmup* for the newly introduced parameters by comparing runs with versus without the re-warmup schedule, while all other components of our framework applied.

Results. Figure 8 shows the conclusions on Muon. In Figure 8(a), RMS-preserved scaling produces a stable and consistent improvement over naive unscaled initialization, ultimately converging to a lower final loss under the same training budget. In Figure 8(b), enabling asymmetric re-warmup further improves late-stage convergence over any other counterparts without re-warmup. Taken together, these results demonstrate that both RMS-preserved scaling and re-warmup remain effective under Muon, confirming that our framework applies beyond AdamW and extends to spectral-style update like Muon without requiring optimizer-specific designs.

H Comparison to Prior Function-Preserving Symmetry-Breaking Heuristics

Experimental setup. Prior width-growth methods that rely on copy-based expansion attempt to break symmetry by two widely used function-preserving heuristics: (i) *Uneven Splitting* [5, 6, 10, 31] by assigning different scaling factors to the channel being copied and the copied one, (ii) *Perturb* [34, 35, 42] by adding symmetric perturbations of equal magnitude and opposite sign to the two duplicated halves. Following the expert-inner expansion in Sec. 4.3, we implement these strategies as well as (iii) *Re-warmup All*, which applies the same re-warmup schedule to all parameters.

Results. Figure 9 shows that these heuristics, despite introducing asymmetry by construction, remain consistently weaker than our framework. Moreover, the zoomed-in view around the expansion moment

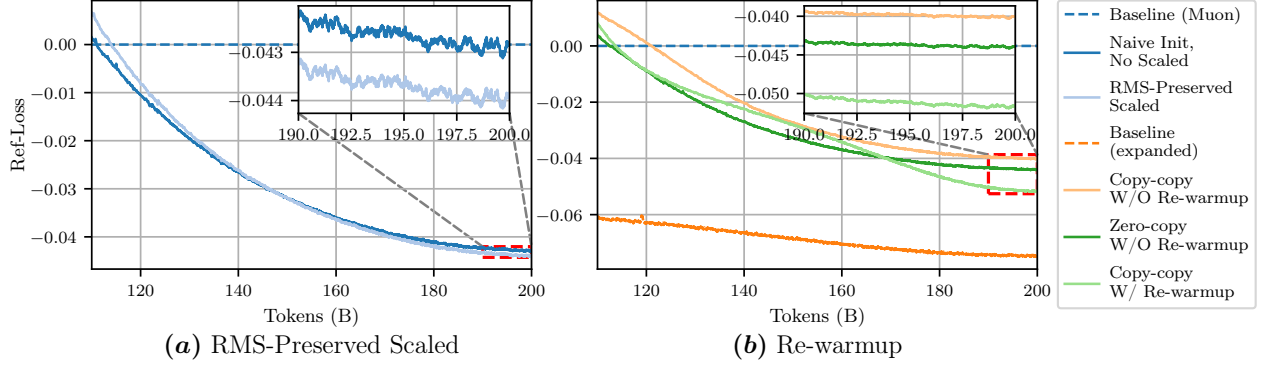


Figure 8 Effectiveness under Muon. We repeat the expert-inner expansion experiment (512→1024 at 100B tokens) using Muon and plot reference-loss versus training tokens. **(a)** RMS-preserved scaling consistently improves late-stage convergence compared to naive unscaled initialization. **(b)** With RMS-preserved scaling and asymmetric state reset applied, asymmetric learning rate re-warmup further lowers the final loss, confirming that our framework remains effective under Muon.

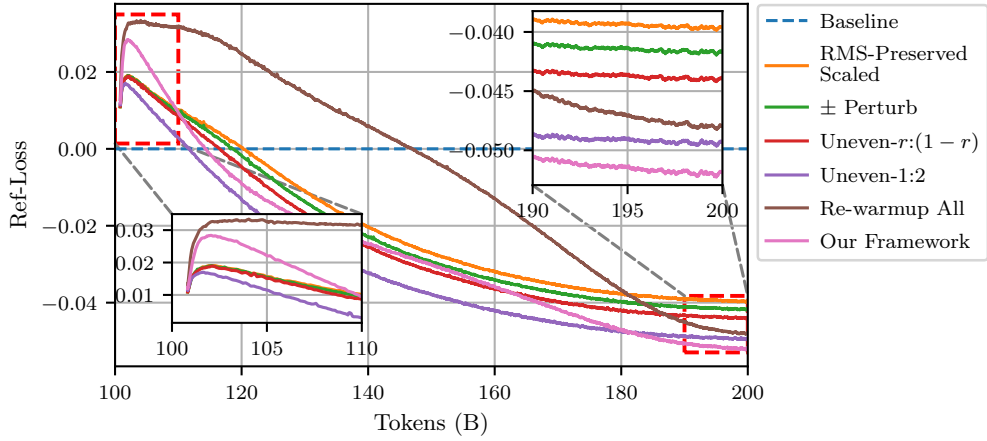


Figure 9 Under expert-inner copy-copy expansion, we compare three alternatives against our framework: (i) Uneven Splitting (fixed 1:2 or randomized $r:(1-r)$ with $r \in [0.1, 0.5]$) and (ii) symmetric \pm perturbation that cancels in the forward pass, and (iii) globally re-warmup all parameters. All alternatives converge to a higher final loss, underperforming our framework. Insets highlight the post-expansion dynamics: our method exhibits a brief loss increase followed by rapid recovery, consistent with more effective symmetry breaking in the newly added capacity.

highlights a transient loss up-shift followed by fast recovery, consistent with targeted exploration for newly introduced parameters benefiting from asymmetric re-warmup.