

U- μ P: THE UNIT-SCALED MAXIMAL UPDATE PARAMETRIZATION

Anonymous authors

Paper under double-blind review

ABSTRACT

The Maximal Update Parametrization (μ P) aims to make the optimal hyperparameters (HPs) of a model independent of its size, allowing them to be swept using a cheap proxy model rather than the full-size target model. We present a new scheme, u- μ P, which improves upon μ P by combining it with Unit Scaling, a method for designing models that makes them easy to train in low-precision. The two techniques have a natural affinity: μ P ensures that the scale of activations is independent of model size, and Unit Scaling ensures that activations, weights and gradients begin training with a scale of one. This synthesis opens the door to a simpler scheme, whose default values are near-optimal. This in turn facilitates a more efficient sweeping strategy, with u- μ P models reaching a loss that is equal to or lower than comparable μ P models and working out-of-the-box in FP8.

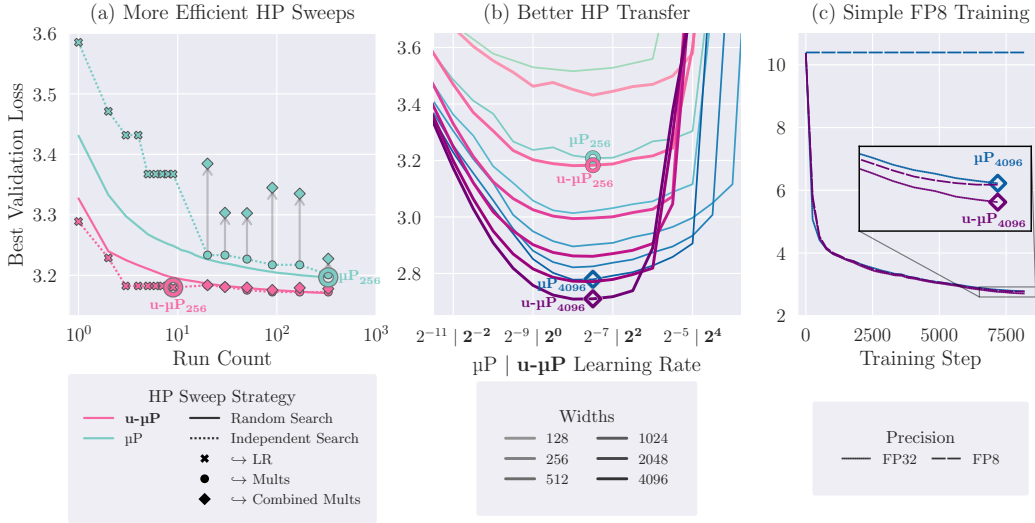


Figure 1: **(a)** Two different HP sweeping processes used for μ P and u- μ P proxy models. Unlike μ P, u- μ P admits independent (1D) search due to careful HP design. The first part of independent search is an LR sweep, which alone reaches near-optimal loss for u- μ P. **(b)** Using the best proxy HPs from (a), we train many models at different widths and LR. The best LR for width 256 is ~optimal for 4096, showing LR transfer along with lower loss. **(c)** We re-train with a simple un-scaled `.to(float8)` cast on matmul inputs. This would fail for other models, but u- μ P trains with minimal degradation.

1 INTRODUCTION

The challenges of large-model training extend beyond the domain of engineering; they are also *algorithmic* in nature. Effective approaches for training smaller models are not guaranteed to work at the multi-billion-parameter scale used for today’s large language models (LLMs). These difficulties can be framed in terms of stability, which we consider in three forms:

1. feature learning stability, which ensures that parts of the model do not learn too fast or slow relative to each other.
2. hyperparameter stability, which ensures that the optimal HPs for small models remain unchanged as the model size grows.
3. numerical stability, which ensures that floating-point representations during training stay within the range of a given number format.

The Maximal Update Parametrization (μP) (Yang & Hu, 2021; Yang et al., 2022) targets the first two sources of instability. μP defines a set of scaling rules that in principle make a model’s optimal HP values consistent across model sizes and ensure ‘maximal feature learning’ in the infinite-width limit. The practical benefits of this are that models continue to improve as they get larger, and that practitioners can re-use a set of HP values (especially the learning rate) found for a small *proxy* version of their model, on a larger *target* model. This is vital for modern LLM training, where the cost of sweeping over candidate HP values for the target model is prohibitive. Consequently, μP has been adopted by several open LLM training efforts (Dey et al., 2023a;b; Liu et al., 2023; Hu et al., 2024) and there are indications of its use in state-of-the-art LLMs¹.

However, there exists a gap between the extensive theory underpinning μP and its effective use in practice. This relates to issues surrounding efficient HP search, HP transfer, interpretability, ease-of-use and low-precision training. Some of these problems have been observed in the literature (Yang et al., 2022; Almazrouei et al., 2023; Lingle, 2024); others we outline here for the first time. As a result, μP does not necessarily provide the kind of simple, stable scaling for which a user might hope.

To address this, we propose the Unit-Scaled Maximal Update Parametrization ($u\text{-}\mu\text{P}$). $u\text{-}\mu\text{P}$ combines μP with another closely-related training innovation, Unit Scaling (Blake et al., 2023). μP ideally provides consistent training dynamics across model sizes, but says little about what those dynamics should be. Unit Scaling addresses this by proposing an ideal principle for dynamics: unit variance for all activations, weights and gradients. Unit Scaling was initially designed to ensure stable numerics, but in the context of μP the principle of unit-scale brings many additional benefits. We show that it provides a foundation upon which the broader range of drawbacks identified for μP can be addressed.

2 BACKGROUND

2.1 THE MAXIMAL UPDATE PARAMETRIZATION

Tensor Programs V (Yang et al., 2022) defines a parametrization as ‘a rule for how to change [HPs] when the widths of a neural network change’. They show that μP is the only parametrization giving ‘maximal feature learning’ in the limit, while standard parametrization (SP) has imbalanced learning.

One consequence of this improved stability is that learning dynamics under μP are ideally independent of model-size, as are optimal HPs. This facilitates a method known as $\mu\text{Transfer}$, which describes the process of training many smaller proxy models to evaluate candidate HP values, then using the best-performing ones to train a larger target model.

ABC-parametrizations. μP , SP, and the Neural Tangent Kernel (NTK) (Jacot et al., 2018) are all instances of abc-parametrizations. This assumes a model under training where weights are defined as:

$$\begin{aligned} W_t &= A_W \cdot w_t, \quad w_0 \sim \mathcal{N}(0, B_W^2), \\ w_{t+1} &= w_t + C_W \cdot \Phi_t(\nabla \mathcal{L}_0, \dots, \nabla \mathcal{L}_t), \end{aligned} \quad (1)$$

with t a time-step and $\Phi_t(\nabla \mathcal{L}_0, \dots, \nabla \mathcal{L}_t)$ is the weight update based on previous loss gradients.

A parametrization such as μP is then defined by specifying how scalars A_W, B_W, C_W change with model width. This can be expressed in terms of width-dependent factors a_W, b_W, c_W , such that $A_W \propto a_W, B_W \propto b_W, C_W \propto c_W$. The values these factors take are what characterize a particular scheme. For μP these are given in Table 1. For depth, a similar result has been proved using depth- μP (Yang et al., 2023b), albeit in a restricted setting. When we refer to μP in the paper we assume the depth- μP scaling rules (Table 2, ‘Residual’ column).

¹ The GPT-4 technical report (OpenAI, 2023) hints at the use of μP by including Yang et al. (2022) in its references, without citing it directly. The multipliers present in the Grok (xAI, 2024) codebase also suggest the use of μP .

A key property of the abc-parametrization is that one can shift scales between A_W, B_W, C_W in a way that preserves learning dynamics (i.e. the activations computed during training are unchanged). We term this *abc-symmetry*. For a fixed $\theta > 0$, the behavior of a network trained with Adam is invariant to changes of the kind:

$$A_W \leftarrow A_W / \theta, \quad B_W \leftarrow B_W \cdot \theta, \quad C_W \leftarrow C_W \cdot \theta \quad (2)$$

(reproduced from Tensor Programs V, Section J.2.1). This means that parametrizations like μP can be presented in different but equivalent ways. ABC-symmetry is a key component in developing u- μP .

Transferable HPs. μP focuses on the subset of HPs whose optimal values we expect to *transfer across* axes such as width and depth. We term these $\mu\text{Transferable HPs}$. All $\mu\text{Transferable HPs}$ function as multipliers and can be split into three kinds: $\alpha_W, \sigma_W, \eta_W$ where $A_W \propto \alpha_W, B_W \propto \sigma_W, C_W \propto \eta_W$. The difference between these multipliers and the ones that define a parametrization is that they are specified by the user, rather than being a function of width. α_W and η_W are rarely introduced outside of the μP literature, but can be valuable to tune for both μP and SP models. In the μP literature the term ‘HPs’ often implicitly refers to $\mu\text{Transferable HPs}$. We adopt this convention here, unless specified otherwise.

Base shape. Two additional non- $\mu\text{Transferable HPs}$ introduced by μP are the base-width and base-depth. This refers to a mechanism where a user specifies a particular shape for the model, where its behavior under μP and SP are the same. The μP model still *scales* according to the abc-rules, so for all other shapes the two models will be different.

In summary, the absolute expressions for A_W, B_W, C_W under μP are given by:

$$A_W \leftarrow \alpha_W \frac{a_W}{a_{W_{\text{base}}}}, \quad B_W \leftarrow \sigma_W \frac{b_W}{b_{W_{\text{base}}}}, \quad C_W \leftarrow \eta_W \frac{c_W}{c_{W_{\text{base}}}} \quad (3)$$

Though base shapes are necessary for μP , they are not typically swept. Rather, they are considered a preference of the user, who may wish to retain the behavior of an existing SP model at a given shape.

Choosing HPs to sweep. In theory, the search space of $\mu\text{Transferable HPs}$ includes $\alpha_W, \sigma_W, \eta_W$ for every parameter tensor W in the model. In practice, far fewer HPs are swept, with global grouping often used for σ_W and η_W , and many α_W are dropped or grouped across layers.

The sets of HPs chosen for sweeps in the μP literature is explored in Appendix E.1. Tensor Programs V uses a random search to identify the best HP values, which has become the standard approach to sweeping. The number of runs in a sweep is typically in the low 100s, incurring a non-negligible cost (though usually less than a single training run of the target model). This high number partly owes to dependencies between HPs (shown in Section 4.1), making the search space hard to explore.

2.2 LOW-PRECISION TRAINING

All the major potential bottlenecks of model training—compute, communication and storage—see roughly linear improvements as the bit-width of their number format is reduced. In modern LLM training, the compute cost of large matrix multiplications (matmuls) means that substantial gains are available if these can be done in low-precision (< 32 bit) formats, which makes them one of the most promising avenues towards increased efficiency in deep learning.

Table 1: The scaling rules defining μP . The type of a weight is determined by whether fan-in & fan-out both depend on width (hidden), only fan-out does (input), or only fan-in (output). Hence fan-in is always a multiple of width here.

ABC-multiplier			Weight (W) Type		
			Input	Hidden	Output
μP	parameter	(a_W)	1	1	$1/\text{fan-in}(W)$
	initialization	(b_W)	1	$1/\sqrt{\text{fan-in}(W)}$	1
	Adam LR	(c_W)	1	$1/\text{fan-in}(W)$	1

Recent AI hardware offers substantial acceleration for the 8-bit FP8 E4 and E5 formats. However their reduced range means that they cannot directly represent some values generated during training. Various methods have been introduced to address this, such as the per-tensor dynamic re-scaling in Transformer Engine (NVIDIA, 2024b). However, this comes at the cost of added complexity and potential overheads. For a more in-depth treatment of low-precision formats, see Appendix J.

2.3 UNIT SCALING

An alternative approach to low-precision training is Unit Scaling (Blake et al., 2023), which also uses per-tensor scaling factors to control range, but instead finds these factors via an analysis of expected tensor statistics at initialization. These are fixed factors, calculated independently of the contents of a tensor, at the beginning of training. As such, the method is easy to use and only adds the overhead of applying static scaling factors (which we show to be negligible in Appendix K).

These factors are chosen to ensure the unit variance of activations, weights and gradients at initialization. This is a useful criterion as it places values around the center of floating-point formats’ absolute range. This applies to all tensors, meaning every operation in the network requires a scaling factor that ensures unit-scaled outputs, assuming unit-scaled inputs. Unit Scaling does not provide a mechanism for re-scaling tensors dynamically during training, but due to its ideal starting scale for gradients, activations and weights this may not be required. Empirically this is shown to be true across multiple architectures, though it is not guaranteed.

We provide an example of deriving the Unit Scaling rule for a matmul op in Appendix E.2, resulting in the scaling factor: $1/\sqrt{d_{\text{fan-in}}}$. We accompany this example with a full recipe for applying Unit Scaling to an arbitrary model.

3 THE UNIT-SCALED MAXIMAL UPDATE PARAMETRIZATION

In this section we show how μP can be adapted to satisfy Unit Scaling, and provide a new set of HPs which—thanks to Unit Scaling—are more interpretable and separable than those commonly used for μP , unlocking several practical benefits. Although some features of u- μP presented here are transformer specific, we stress that u- μP can in principle be applied to a wide range of architectures, since both μP and Unit Scaling are very general approaches.

3.1 COMBINING μP WITH UNIT SCALING

Whereas Unit Scaling provides rules for scaling all operations, μP only does so for parametrized ones. It’s these operations we need to address to arrive at a unified scheme, resolving differences in the scaling rules each recommends. We begin with the expressions for the A_W, B_W, C_W scaling factors in Equation (3), and substitute in the μP scaling rules defined in Table 1. This results in a full implementation of μP , which is shown in the top half of Table 2. We set out to turn this into a valid Unit Scaling scheme, which requires unit initializations ($B_W \leftarrow 1$) and matmuls with the Unit Scaling factor we identified in Section 2.3 ($A_W \leftarrow 1/\sqrt{\text{fan-in}}$).

Our first step is to drop the σ_W and base-fan-in HPs entirely, and associate the α_W HPs with subsequent non-linear functions instead of weights—decisions we justify in the rest of this section (this results in the simplified intermediate implementation in Table 11). Our input weights now have unit initializations as desired, and a unit parameter multiplier, which is also the appropriate scaling factor (as input layers here are embedding lookups, not matmuls).

Hidden weights now have the implementation: $A_W \leftarrow 1, B_W \leftarrow \frac{1}{\sqrt{\text{fan-in}}}, C_W \leftarrow \eta \frac{1}{\text{fan-in}}$, which differs from our Unit Scaling criteria. However, using abc-symmetry (Equation (2)) we can shift scales by $\sqrt{\text{fan-in}}$, arriving at a unit-scaled scheme: $A_W \leftarrow \frac{1}{\sqrt{\text{fan-in}}}, B_W \leftarrow 1, C_W \leftarrow \eta \frac{1}{\sqrt{\text{fan-in}}}$.

Finally, our output layers also have unit initialization, but a parameter multiplier of $A_W \leftarrow 1/\text{fan-in}$. This differs from the Unit Scaling rule, but in the forward pass this is permissible as there are no subsequent matmuls in a transformer. In the backward pass this mis-scaling would propagate, so we apply the desired $\leftarrow 1/\sqrt{\text{fan-in}}$ factor. Using different forward and backward scales in this way is usually not allowed, but is valid for output layers due to the cut-edge rule (Appendix H).

Table 2: Scaling rules for μP versus $\text{u-}\mu\text{P}$, including associated HPs (assuming the *extended* set in Table 3). These rules constitute the definition of $\text{u-}\mu\text{P}$, along with the unit-scaled ops in Appendix B.

	ABC-multiplier	Weight Type			Residual
		Input	Hidden	Output	
μP	parameter (A_W)	α_{emb}	1 (or α_{attn})	$\alpha_{\text{out}} \frac{\text{base-fan-in}}{\text{fan-in}}$	$\sqrt{\frac{\text{base-depth}}{\text{depth}}}^*$
	initialization (B_W)	σ_{init}	$\sigma_{\text{init}} \sqrt{\frac{\text{base-fan-in}}{\text{fan-in}}}$	σ_{init}	—
	Adam LR (C_W)	$\eta \hat{\eta}_{\text{emb}}$	$\eta \frac{\text{base-fan-in}}{\text{fan-in}}$	η	$\sqrt{\frac{\text{base-depth}}{\text{depth}}}$
$\text{u-}\mu\text{P}$	parameter [†] (A_W)	1	$\frac{1}{\sqrt{\text{fan-in}}}$	$\frac{1}{\text{fan-in}}^{\ddagger}$	$\frac{1}{\sqrt{\text{depth}}}^*$
	initialization (B_W)	1	1	1	—
	Adam LR (C_W)	$\eta \frac{1}{\sqrt{\text{fan-out}}}$	$\eta \frac{1}{\sqrt{\text{fan-in}}}$	η	$\frac{1}{\sqrt{\text{depth}}}$

*Residual multipliers are applied to the end of each branch, rather than the output of linear layers.

[†] $\text{u-}\mu\text{P}$'s α HPs are associated with operations, not weights, so are not included here (see Section 3.3).

[‡]To maintain unit scale we apply $1/\sqrt{\text{fan-in}}$ scaling in the backward pass (see Appendix H).

The final change we make is to the input LR scaling rule, which we show in Section 3.4 is more effective if $c_W \leftarrow 1$ is replaced with $c_W \leftarrow 1/\sqrt{\text{fan-out}}$. With these changes made, we arrive at our final $\text{u-}\mu\text{P}$ scheme, given in Table 2. Note that the scaling rules in this table must be combined with the standard Unit Scaling rules for other non-matmul operations. These are covered in Appendix B.

3.2 OUT-OF-THE-BOX LOW-PRECISION TRAINING

When training a transformer model with $\text{u-}\mu\text{P}$ most tensors have stable scale during training, except a small number of *critical tensors* that exhibit scale growth, which can cause extreme values to go out of FP8 range. We empirically identify these to be the inputs to the attention dense projection and final FFN matmul as well as the weight of the decoder² (for details see Appendix A.7).

Based on our analysis, we propose the following FP8 mixed-precision scheme for $\text{u-}\mu\text{P}$ transformers:

- For layers using non-critical tensors, we cast the input and weight to E4M3, and the gradient with respect to the output to E5M2. This is done in the forward computation, as well as the two backward computations (for the gradients w.r.t. the weight and input).
- For layers using critical tensors, matmuls are performed in BF16 (along with non-matmul operations). We keep optimizer states in FP32, leaving the use of FP8 to future work.

In some cases we are able to deal with the critical tensors by casting them to E5M2. However, we observed instabilities applying this in a large-scale setting. In small-scale settings we also empirically find that applying E4M3 instead of E5M2 for the gradients is possible, but again becomes problematic in the large-scale setting where gradients require higher dynamic range. Under our mixed-precision scheme, approximately 70% of the matmul computations in a Llama transformer block are performed in FP8. If desired, a dynamic per-tensor scaling could still be applied to the critical tensors.

3.3 A PRINCIPLED APPROACH TO HYPERPARAMETERS

Approaches for selecting which HPs to sweep are poorly motivated in the literature (see Appendix C.2). Our objective in $\text{u-}\mu\text{P}$ is to find a simple, well-justified and effective alternative. To this end, we propose the following ideal criteria:

1. **Minimal cardinality:** the use of as few HPs as possible.

² The decoder becomes negligible in terms of model flops as width and depth of the model increase, so we generally keep this operation in higher precision.

2. **Maximal expressivity**: the ability to still express any model defined using the per-tensor $\alpha_W, \sigma_W, \eta_W$ HPs outlined in Section 2.1 (in practice, we relax this slightly).
3. **Minimal interdependency**: the optimal value of each HP should not depend on the value of other HPs, simplifying the search space.
4. **Interpretability**: there should be a clear explanation for what an HP’s value ‘means’ in the context of the model.

The u- μ P HPs given in Table 3 are designed to satisfy these criteria, to the fullest extent possible. The placement of these HPs in the model is given in Table 7.

Cardinality & expressivity. We arrive at our set of HPs in three steps, starting with the full $\alpha_W, \sigma_W, \eta_W$ for each weight tensor W . Firstly, we can remove any one of these HPs by permuting under abc-symmetry, such that one HP = 1. As we want our weights to begin with unit scale, we choose to drop σ_W . Secondly, we observe that several of the α_W HPs combine linearly with other α_W HPs, providing an opportunity to re-parametrize with a single HP. For instance, the scale of self-attention softmax activations is proportional to the product of α_W multipliers: $\text{std}(x_{\text{attn}}) \propto \alpha_{W_Q} \alpha_{W_K}$. In this instance we use a single α parameter (termed $\alpha_{\text{attn-softmax}}$) and associate it with the attention operation, rather than the weights.

We apply the same principle to all operations, unless they are unary and k -homogeneous for $k \geq 0$, in which case they propagate scale and don’t require an HP (see Appendix G.1). This results in the set of HPs shown, with their placement in the model given in Table 7. We note that this procedure also works in other architectures than transformers, and naturally will produce other kinds of multipliers.

Thirdly, we use a single global η and group α HPs across layers. This breaks our expressivity criterion, but we argue represents the best trade-off between expressivity and cardinality. We show in Appendix A.3 that having tuned a global η HP and our extended α HPs, the further benefits of tuning per-tensor $\hat{\eta}_W$ HPs (which modify the global η) is minimal.

Interdependency. The second stage above, moving α HPs from weights into operations, not only reduces the number of HPs, but also minimizes the interdependence between those that remain. We find that u- μ P’s optimal HP values depend less on each other than under μ P (see Section 4.1).

Interpretability. The combination of unit scale and reduced dependencies between HPs means that each α can be interpreted as determining some fundamental property of the model at initialization. For example, the $\alpha_{\text{loss-softmax}}$ HP defines the (inverse of) the softmax’s *temperature* for a unit-scaled input. We also introduce a new scaling scheme (defined in Appendix G.2.2) for residual connections, designed to give HPs independence and a clear interpretation: α_{res} defines the contribution of the residual connections to the output scale, and $\alpha_{\text{res-attn-ratio}}$ defines the relative contribution of attention versus FFN branches. Finally, we choose not to include base shape HPs in u- μ P. They do not add to expressivity, lack a clear interpretation (besides alignment to a base model at a particular shape), break the interpretation of α HPs given above, and complicate implementation.

3.4 A NEW EMBEDDING LR RULE

Although theoretical transfer properties have been proved for μ P, not all its HPs have had μ Transfer shown empirically. We do so for the *extended* μ P transformer HPs in Figure 15, where we observe poor transfer across width for the embedding LR multiplier $\hat{\eta}_{\text{emb}}$. This suggests that the corresponding LR scaling rule for c_{emb} is mis-specified.

We show in Figure 2 that changing it from the μ P rule of $c_{\text{emb}} = 1$ to $1/\sqrt{\text{fan-out}}$ corrects this failure of HP transfer. As a result, we see improved loss under u- μ P for larger model-sizes relative to μ P. Our adoption of this change is a key factor in the improved performance of u- μ P over μ P in Figure 1. We offer no theoretical justification for this change, which we leave to future work.

Table 3: Typical transformer HPs used under different schemes. *Basic* HPs in **bold** are most impactful and are commonly swept. *Extended* HPs in non-bold are not always swept, often set heuristically or dropped.

SP	μ P	u- μ P
η	η	η
σ -scheme	σ_{init}	
	$\alpha_{\text{emb}} \eta_{\text{emb}}$	$\alpha_{\text{ffn-act}}$
	α_{attn}	$\alpha_{\text{attn-softmax}}$
	α_{out}	α_{res}
	base-width	$\alpha_{\text{res-attn-ratio}}$
	base-depth	$\alpha_{\text{loss-softmax}}$

3.5 HYPERPARAMETER SEARCH

As shown in section Section 2.1, the standard approach to HP search for μ Transfer is via a random sweep over all HPs simultaneously. Sweeping individual HPs separately is challenging due to the dependencies between them. In contrast, u- μ P’s HPs are designed to admit such a strategy due to our interdependence criterion. Because of this, we propose a simpler sweeping strategy for u- μ P which we term *independent search* (outlined in detail in Appendix A.5). Independent search involves a sweep of the LR, followed by a set of one-dimensional sweeps of the other HPs (which can be run in parallel). The best results from the individual sweeps are combined to form the final set of HP values.

4 EXPERIMENTS

Our experiments use the Llama architecture (Touvron et al., 2023a) trained on WikiText-103 (Merity et al., 2017) (except large-scale runs in Section 4.4). We apply best-practice LLM training techniques from the literature (see Table 5). In accordance with our analysis in Appendix C.1, we remove parameters from norm layers, use independent AdamW, and avoid training on too many epochs.

4.1 QUANTIFYING HYPERPARAMETER INTERDEPENDENCE

Our principled approach to HPs (see Section 3.3) contains the requirement that their optimal values should depend minimally on the value of other HPs. We now investigate this empirically, conducting a 2D sweep over every pair of HPs for μ P and u- μ P (see Figures 11 and 12 for pairwise results).

To derive an empirical measure of HP dependency, we introduce the notion of *transfer error* (see Algorithm 1). We take the best value of the transfer HP for each non-optimal value of the fixed HP, and use it with the optimal value of the fixed HP. The transfer error is the difference between the losses obtained and the minimum loss. Figure 3 shows this measure for each pair of HPs, reflecting the improvement in HP dependency as a result of our scheme. This gives u- μ P a reduced risk of small transfer errors leading to large degradations, and the potential to sweep HPs in a more separable way.

4.2 HYPERPARAMETER SEARCH

Our new HP scheme, designed for improved separability, enables proxy model HPs to be swept more efficiently. This is shown in Figure 1 (a). We conduct a standard random search for μ P and u- μ P, along with the independent search outlined in Section 3.5 (and Appendix A.5).

Independent search begins with a simple LR sweep. This alone is sufficient for u- μ P to reach near-optimal loss (using only 9 runs). During this phase other HPs are fixed at 1, which for u- μ P means that the inputs to operations are generally unit-scaled. Consequently, we conclude that unit scale at

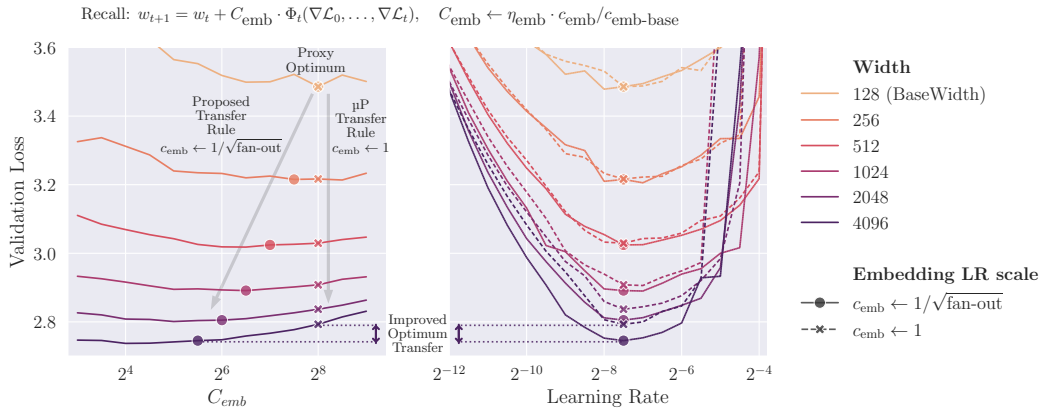


Figure 2: (Left) holding the embedding LR multiplier (C_{emb}) constant, vs. scaling with $\sqrt{1/\text{width}}$, both with a fixed global LR. This suggests the μ P embedding LR rule (c_{emb}) should follow the latter scaling. (Right) we test this by sweeping the global LR under the two scaling rules. The new rule leads to lower loss on large models. (Dot/cross markers represent the same runs across both graphs).

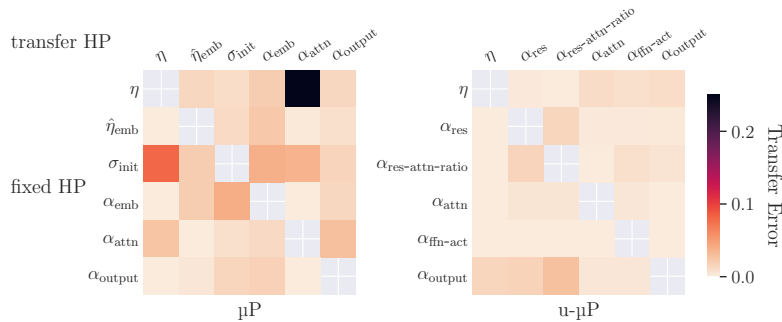


Figure 3: A visualization of the dependencies between pairs of HPs under each scheme. Transfer error

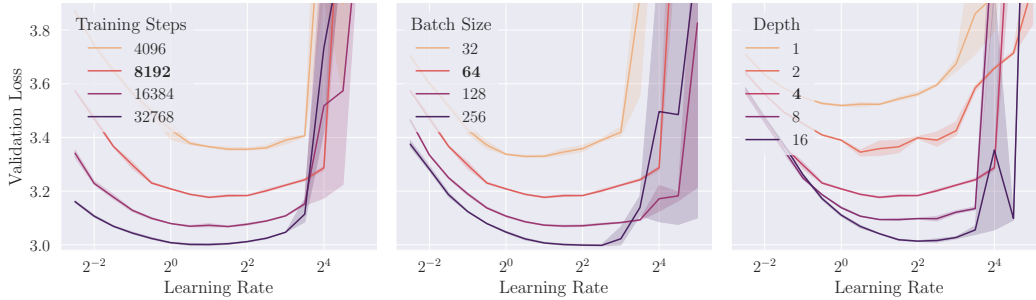


Figure 4: Learning rate transfer for u- μ P over training steps, batch size and depth. See Figure 1 (b) for transfer over width and Figure 13 for regular μ P transfer results. The **default** shape parameter for other panels is shown in bold. The shaded area shows the 95% confidence interval for the mean.

initialization is close to ideal scaling for effective learning here. In contrast μP still requires non-LR HPs to be swept to attain a reasonable loss. The final ‘combined mults’ phase causes the loss to spike for μP . This is due to the HP dependencies shown in Figure 3, which mean HPs cannot be swept independently and used together, necessitating random search which can require hundreds of runs.

4.3 HYPERPARAMETER TRANSFER

We train many models and plot transfer of LR across width (Figure 1 (b)), steps, batch size and depth (Figure 4), and transfer of other HPs across width (Figure 15). Note that u- μ P (building on μ P) is designed to give transfer over width³; the other axes we report for practical purposes. We find that:

1. The optimal LR is constant across width under u- μ P. There is a small drift for training steps and batch size, and a larger one with depth. Hence we recommend proxy models which primarily differ in width, moderately in steps and batch size, and least in depth.
2. Whereas μ P sees diminishing returns for larger widths, u- μ P continues to benefit from width. We attribute this primarily to our improved embedding LR rule.
3. Non-LR HPs also have constant optima across width under u- μ P. This is not true for μ P, where $\hat{\eta}_{\text{emb}}$ has poor transfer (see Section 3.4), along with σ_{init} (see Appendix C.2).
4. The optimal values found for non-LR HPs are all close to 1. In practice this means that dropping these HPs entirely is potentially viable for similar models and training setups.

4.4 FP8 TRAINING

In this section we justify the simple mixed-precision scheme described in Section 3.2 and demonstrate that it can be used to train u-uP models out-of-the-box.

Proof-of-concept. Figure 5 shows the RMS of all linear layer inputs for a moderately sized transformer. RMS captures the larger of the mean and scale of a distribution, and as such is a good test of whether a tensor is likely to suffer over/underflow in low-precision. We observe that u- μ P tensors largely have RMS starting close to 1 and remaining so at the end of training, supporting our scheme.

³ As we use depth- μP this could be said about depth as well, but as Yang et al. (2023b) show that transformers don’t attain depth-transfer under depth- μP we do not expect strong transfer across depth.

Table 4: 0-shot benchmark results at 7B scale.

Scheme	Format	MMLU	HellaSwag	OpenBook QA	PIQA	TriviaQA	WinoGr
SP	BF16	29.6	52.4	27.8	76.5	22.2	63.3
u- μ P	BF16	29.0	53.4	31.6	77.1	23.4	63.7
u- μ P	FP8	31.2	53.4	29.6	77.6	21.3	65.7

Figure 17 demonstrates the scale-growth of critical tensors which our scheme is designed to accommodate, showing RMS on a per-tensor basis over steps. Figure 18 provides further insight into this issue, showing the effect of LR, width, depth, steps and batch size on the RMS of critical tensors.

As an initial proof-of-concept we train a u- μ P model using our FP8 scheme over 8k steps, using HPs from a proxy model, as shown in Figure 1 (c). We see only a small degradation versus FP32, and at this scale critical tensors can still be cast to FP8 using E5M2, while gradients can even use E4M3.

Larger scale. Next we consider a more realistic training scenario. Using the same architecture, and following the steps set out in our u- μ P user-guide (Appendix D), we train our target models on 300B tokens of the SlimPajama dataset (Shen et al., 2023) (see Appendix A.8 for training details).

We begin with an independent search (Section 3.5) over our u- μ P proxy model’s HPs. Here we make the following observations:

1. When using a relatively small proxy model (8 layers and 512 width), the HP-loss landscape is rather noisy. By doubling the width we can discern optimal HP values more clearly.
2. The most important HPs are η and $\alpha_{\text{res-attn-ratio}}$. All others can be left at the default of 1.
3. The optimal values of these HPs are $\eta = 2^{3.5}$ and $\alpha_{\text{res-attn-ratio}} = 2^{-2.0}$ and thus differ non-trivially from the observed HPs in our smaller-scale experiments.

We then train u- μ P models of approximately 1B, 3B and 7B parameters, using our FP8 mixed-precision scheme (see Section 3.2). We also train two baselines at each size: the first is a BF16 version of our u- μ P models, and the second is a set of SP models using the weight init scheme from the Pythia model family (Biderman et al., 2023) and the LR scheme from Llama 3 (Dubey et al., 2024), scaling inversely with width and using a LR of 3e-4 at 7B scale. For the SP baseline, we use a non-independent weight decay of 0.1 and parametric RMS norm as in Llama (Dubey et al., 2024), which is standard practice. The loss curves are shown in Figure 6. All FP8 runs converge and show no significant loss degradation. In comparison to SP, the u- μ P models have a qualitatively different training curve with a higher loss for most of training that catches up in latter stages, hinting at a fundamentally different optimization trajectory. In terms of downstream performance, both of the u- μ P 7B models are competitive with SP. In particular, the scores of the FP8 model are mostly on par with the BF16 models (see Table 4).

5 RELATED WORK

Low-precision training Techniques to facilitate FP8 training include those covered in Appendix J and more (Wang et al., 2018; Mellempudi et al., 2019; Perez et al., 2023). These largely concern the quantizing of activations, weights and gradients, though Peng et al. (2023) also explore FP8 optimizer

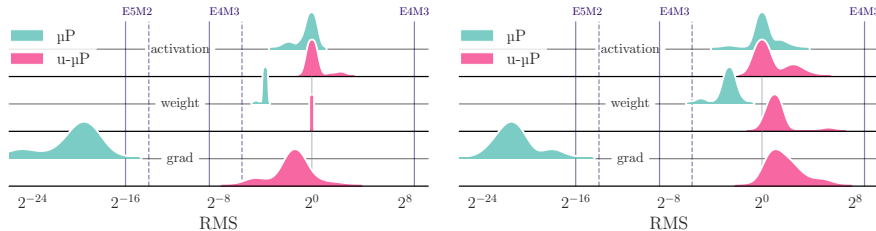


Figure 5: Per-tensor $\text{RMS} = \sqrt{\sigma^2 + \mu^2}$ across u- μ P and μ P models at initialization (left) and after training (right). Dashed and solid red lines show each format’s min. normal and subnormal values.

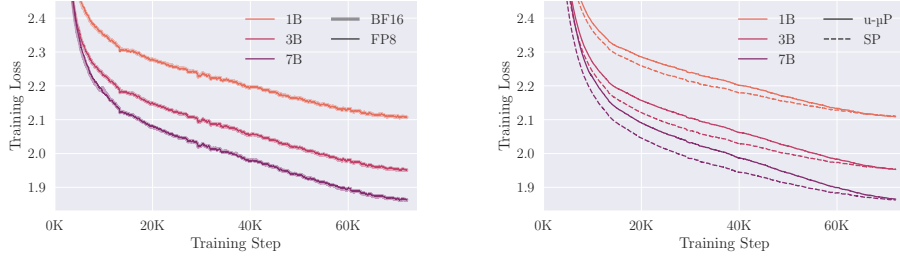


Figure 6: Large-scale training runs. (Left) u-μP BF16 vs u-μP FP8. (Right) u-μP BF16 vs SP BF16.

states and cross-device comms, which we consider interesting avenues of exploration. Quantization is made harder by the emergence of large outlier values when training at scale, with several techniques proposed to mitigate this (Bondarenko et al., 2023; Sun et al., 2024; He et al., 2024). Combining these with u-μP is a natural extension of our work, facilitating simpler or lower-bit quantization.

Parametrizations The neural tangent kernel (NTK) (Jacot et al., 2018) and mean field parametrization (MFP) (Mei et al., 2018; Bordelon & Pehlevan, 2023) are alternatives to μP. In terms of width exponents, u-μP is identical to MFP except for the embedding learning rate scaling, as shown in Table 1 of Everett et al. (2024). The equivalence classes given by Everett et al. (2024) also mirror our notion of abc-symmetry. They additionally show that under MFP gradients become sufficiently small to require re-scaling the Adam epsilon term. u-μP naturally solves this problem in a similar way, via the up-scaling of our gradients to unit scale. Another difference between MFP and u-μP is that the latter specifies absolute A_W, B_W scales rather than just exponents, which is key to numerical stability.

Compatible parametrization frameworks Apart from u-μP, several other recent efforts have also introduced frameworks for training under particular parametrizations, which have the potential to be compatible with u-μP. Large et al. (2024) introduces the *modular norm* over the full weight-space, which like μP aims to ensure stable updates that provide LR transfer, and like u-μP is implemented via modules designed for stable training. Everett et al. (2024) explore the notion of *alignment* between parameters and data, showing that other parametrizations with per-layer LRs can outperform standard μP. These parametrizations could admit unit scaling under abc-symmetry, though this is outside the scope of this work. Recent work showing how to apply μP in the settings of weight sparsity Dey et al. (2024a) and structured matrices Qiu et al. (2024) are also interesting candidates for u-μP.

Signal propagation Unit Scaling and μP are mainly concerned about statistical properties of single tensors. Another line of work studies the covariance between two distinct inputs and how it propagates through the model. For MLPs this was analyzed in Poole et al. (2016); Schoenholz et al. (2017), revealing phase transitions in deep networks related to expressivity and trainability. For transformers, more recent work (Noci et al., 2022) studies the collapse of token representations. Coming from this angle, they derive the same $1/\sqrt{\text{depth}}$ residual multiplier scaling that is used in u-μP and identify the attention softmax temperature as a HP to control the gradient scale of query and key weights.

6 CONCLUSIONS

We introduce u-μP, a modified and improved version of μP that satisfies Unit Scaling. Through careful analysis guided by first principles we identify an interpretable set of HPs that has minimal interdependencies and facilitates an efficient independent sweeping strategy. We show that the stability properties of μP combined with Unit Scaling enable a simple and robust FP8 mixed precision scheme that works in a realistic large scale training scenario.

Limitations and future work. Some choices like the modified embedding LR rule are only justified by empirical observations, and currently lack a theoretical explanation. Additionally, neither μP nor Unit Scaling give guarantees for network quantities to be well-behaved over the course of training. In particular we would like to understand the issue (or feature) of scale growth in the critical layers better and look into possible mitigations. We also believe that low-precision training techniques can be pushed further, with u-μP offering an ideal starting point for future optimizations.

REFERENCES

- Ebtesam Almazrouei, Hamza Alobeidli, Abdulaziz Alshamsi, Alessandro Cappelli, Ruxandra Cojocaru, Mérouane Debbah, Étienne Goffinet, Daniel Hesslow, Julien Launay, Quentin Malartic, Daniele Mazzotta, Badreddine Noune, Baptiste Pannier, and Guilherme Penedo. The Falcon series of open language models. *CoRR*, abs/2311.16867, 2023. doi: 10.48550/ARXIV.2311.16867. URL <https://doi.org/10.48550/arXiv.2311.16867>.
- Anonymous. Straight to zero: Why linearly decaying the learning rate to zero works best for LLMs. In *Submitted to The Thirteenth International Conference on Learning Representations*, 2024. URL <https://openreview.net/forum?id=hrOlBgHsMI>. under review.
- Stella Biderman, Hailey Schoelkopf, Quentin Anthony, Herbie Bradley, Kyle O’Brien, Eric Hallahan, Mohammad Aflah Khan, Shivanshu Purohit, USVSN Sai Prashanth, Edward Raff, Aviya Skowron, Lintang Sutawika, and Oskar van der Wal. Pythia: A suite for analyzing large language models across training and scaling, 2023. URL <https://arxiv.org/abs/2304.01373>.
- Charlie Blake, Douglas Orr, and Carlo Luschi. Unit scaling: Out-of-the-box low-precision training. In Andreas Krause, Emma Brunskill, Kyunghyun Cho, Barbara Engelhardt, Sivan Sabato, and Jonathan Scarlett (eds.), *International Conference on Machine Learning, ICML 2023, 23-29 July 2023, Honolulu, Hawaii, USA*, volume 202 of *Proceedings of Machine Learning Research*, pp. 2548–2576. PMLR, 2023. URL <https://proceedings.mlr.press/v202/blake23a.html>.
- Yelysei Bondarenko, Markus Nagel, and Tijmen Blankevoort. Quantizable transformers: Removing outliers by helping attention heads do nothing. In Alice Oh, Tristan Naumann, Amir Globerson, Kate Saenko, Moritz Hardt, and Sergey Levine (eds.), *Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023*, 2023. URL http://papers.nips.cc/paper_files/paper/2023/hash/edbc7583fd8921dad78adecfe06a99b-Abstract-Conference.html.
- Blake Bordelon and Cengiz Pehlevan. Self-consistent dynamical field theory of kernel evolution in wide neural networks*. *Journal of Statistical Mechanics: Theory and Experiment*, 2023 (11):114009, nov 2023. doi: 10.1088/1742-5468/ad01b0. URL <https://dx.doi.org/10.1088/1742-5468/ad01b0>.
- Tri Dao, Daniel Y. Fu, Stefano Ermon, Atri Rudra, and Christopher Ré. Flashattention: Fast and memory-efficient exact attention with io-awareness. In Sanmi Koyejo, S. Mohamed, A. Agarwal, Danielle Belgrave, K. Cho, and A. Oh (eds.), *Advances in Neural Information Processing Systems 35: Annual Conference on Neural Information Processing Systems 2022, NeurIPS 2022, New Orleans, LA, USA, November 28 - December 9, 2022*, 2022. URL http://papers.nips.cc/paper_files/paper/2022/hash/67d57c32e20fd0a7a302cb81d36e40d5-Abstract-Conference.html.
- Tim Dettmers, Mike Lewis, Younes Belkada, and Luke Zettlemoyer. Gpt3.int8(): 8-bit matrix multiplication for transformers at scale. In Sanmi Koyejo, S. Mohamed, A. Agarwal, Danielle Belgrave, K. Cho, and A. Oh (eds.), *Advances in Neural Information Processing Systems 35: Annual Conference on Neural Information Processing Systems 2022, NeurIPS 2022, New Orleans, LA, USA, November 28 - December 9, 2022*, 2022. URL http://papers.nips.cc/paper_files/paper/2022/hash/c3ba4962c05c49636d4c6206a97e9c8a-Abstract-Conference.html.
- Nolan Dey, Gurpreet Gosal, Zhiming Chen, Hemant Khachane, William Marshall, Ribhu Pathria, Marvin Tom, and Joel Hestness. Cerebras-GPT: Open compute-optimal language models trained on the cerebras wafer-scale cluster. *CoRR*, abs/2304.03208, 2023a. doi: 10.48550/ARXIV.2304.03208. URL <https://doi.org/10.48550/arXiv.2304.03208>.
- Nolan Dey, Daria Soboleva, Faisal Al-Khateeb, Bowen Yang, Ribhu Pathria, Hemant Khachane, Shaheer Muhammad, Zhiming Chen, Robert Myers, Jacob Robert Steeves, Natalia Vassilieva, Marvin Tom, and Joel Hestness. BTLN-3B-8K: 7B parameter performance in a 3B parameter model. *CoRR*, abs/2309.11568, 2023b. doi: 10.48550/ARXIV.2309.11568. URL <https://doi.org/10.48550/arXiv.2309.11568>.

- Nolan Dey, Shane Bergsma, and Joel Hestness. Sparse maximal update parameterization: A holistic approach to sparse training dynamics. *CoRR*, abs/2405.15743, 2024a. doi: 10.48550/ARXIV.2405.15743. URL <https://doi.org/10.48550/arXiv.2405.15743>.
- Nolan Dey, Shane Bergsma, and Joel Hestness. Sparse maximal update parameterization: A holistic approach to sparse training dynamics. *CoRR*, abs/2405.15743, 2024b. URL <http://arxiv.org/abs/2405.15743>.
- Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, Anirudh Goyal, Anthony Hartshorn, Aobo Yang, Archi Mitra, Archie Sravankumar, Artem Korenev, Arthur Hinsvark, Arun Rao, Aston Zhang, Aurelien Rodriguez, Austen Gregerson, Ava Spataru, Baptiste Roziere, Bethany Biron, Binh Tang, Bobbie Chern, Charlotte Caucheteux, Chaya Nayak, Chloe Bi, Chris Marra, Chris McConnell, Christian Keller, Christophe Touret, Chunyang Wu, Corinne Wong, Cristian Canton Ferrer, Cyrus Nikolaidis, Damien Allonsius, Daniel Song, Danielle Pintz, Danny Livshits, David Esiobu, Dhruv Choudhary, Dhruv Mahajan, Diego Garcia-Olano, Diego Perino, Dieuwke Hupkes, Egor Lakomkin, Ehab AlBadawy, Elina Lobanova, Emily Dinan, Eric Michael Smith, Filip Radenovic, and Frank Zhang et al. The llama 3 herd of models, 2024. URL <https://arxiv.org/abs/2407.21783>.
- Katie Everett, Lechao Xiao, Mitchell Wortsman, Alexander A. Alemi, Roman Novak, Peter J. Liu, Izzeddin Gur, Jascha Sohl-Dickstein, Leslie Pack Kaelbling, Jaehoon Lee, and Jeffrey Pennington. Scaling exponents across parameterizations and optimizers, 2024. URL <https://arxiv.org/abs/2407.05872>.
- Dirk Groeneveld, Iz Beltagy, Evan Pete Walsh, Akshita Bhagia, Rodney Kinney, Oyvind Tafjord, Ananya Harsh Jha, Hamish Ivison, Ian Magnusson, Yizhong Wang, Shane Arora, David Atkinson, Russell Authur, Khyathi Raghavi Chandu, Arman Cohan, Jennifer Dumas, Yanai Elazar, Yuling Gu, Jack Hessel, Tushar Khot, William Merrill, Jacob Morrison, Niklas Muennighoff, Aakanksha Naik, Crystal Nam, Matthew E. Peters, Valentina Pyatkin, Abhilasha Ravichander, Dustin Schwenk, Saurabh Shah, Will Smith, Emma Strubell, Nishant Subramani, Mitchell Wortsman, Pradeep Dasigi, Nathan Lambert, Kyle Richardson, Luke Zettlemoyer, Jesse Dodge, Kyle Lo, Luca Soldaini, Noah A. Smith, and Hannaneh Hajishirzi. Olmo: Accelerating the science of language models. In Lun-Wei Ku, Andre Martins, and Vivek Srikumar (eds.), *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2024, Bangkok, Thailand, August 11-16, 2024*, pp. 15789–15809. Association for Computational Linguistics, 2024. doi: 10.18653/V1/2024.ACL-LONG.841. URL <https://doi.org/10.18653/v1/2024.acl-long.841>.
- Alexander Hägele, Elie Bakouch, Atli Kosson, Loubna Ben Allal, Leandro von Werra, and Martin Jaggi. Scaling laws and compute-optimal training beyond fixed training durations. *CoRR*, abs/2405.18392, 2024. doi: 10.48550/ARXIV.2405.18392. URL <https://doi.org/10.48550/arXiv.2405.18392>.
- Bobby He, Lorenzo Noci, Daniele Paliotta, Imanol Schlag, and Thomas Hofmann. Understanding and minimising outlier features in neural network training. *CoRR*, abs/2405.19279, 2024. doi: 10.48550/ARXIV.2405.19279. URL <https://doi.org/10.48550/arXiv.2405.19279>.
- Shengding Hu, Yuge Tu, Xu Han, Chaoqun He, Ganqu Cui, Xiang Long, Zhi Zheng, Yewei Fang, Yuxiang Huang, Weilin Zhao, Xinrong Zhang, Zhen Leng Thai, Kai Zhang, Chongyi Wang, Yuan Yao, Chenyang Zhao, Jie Zhou, Jie Cai, Zhongwu Zhai, Ning Ding, Chao Jia, Guoyang Zeng, Dahai Li, Zhiyuan Liu, and Maosong Sun. MiniCPM: Unveiling the potential of small language models with scalable training strategies. *CoRR*, abs/2404.06395, 2024. doi: 10.48550/ARXIV.2404.06395. URL <https://doi.org/10.48550/arXiv.2404.06395>.
- IEEE Computer Society. IEEE standard for floating-point arithmetic. pp. 1–84, July 2019. doi: 10.1109/IEEESTD.2019.8766229.
- Arthur Jacot, Clément Hongler, and Franck Gabriel. Neural tangent kernel: Convergence and generalization in neural networks. In Samy Bengio, Hanna M. Wallach, Hugo Larochelle, Kristen Grauman, Nicolò Cesa-Bianchi, and Roman Garnett (eds.), *Advances in Neural Information Processing Systems 31: Annual Conference on Neural Information*

- Processing Systems 2018, NeurIPS 2018, December 3-8, 2018, Montréal, Canada*, pp. 8580–8589, 2018. URL <https://proceedings.neurips.cc/paper/2018/hash/5a4belfa34e62bb8a6ec6b91d2462f5a-Abstract.html>.
- Oleksii Kuchaiev, Boris Ginsburg, Igor Gitman, Vitaly Lavrukhin, Carl Case, and Paulius Micikevicius. OpenSeq2Seq: Extensible toolkit for distributed and mixed precision training of sequence-to-sequence models. *CoRR*, abs/1805.10387, 2018. URL <http://arxiv.org/abs/1805.10387>.
- Tim Large, Yang Liu, Minyoung Huh, Hyojin Bahng, Phillip Isola, and Jeremy Bernstein. Scalable optimization in the modular norm. *CoRR*, abs/2405.14813, 2024. doi: 10.48550/ARXIV.2405.14813. URL <https://doi.org/10.48550/arXiv.2405.14813>.
- Lucas D. Lingle. A large-scale exploration of μ -transfer. *CoRR*, abs/2404.05728, 2024. doi: 10.48550/ARXIV.2404.05728. URL <https://doi.org/10.48550/arXiv.2404.05728>.
- Zhengzhong Liu, Aurick Qiao, Willie Neiswanger, Hongyi Wang, Bowen Tan, Tianhua Tao, Junbo Li, Yuqi Wang, Suqi Sun, Omkar Pangarkar, Richard Fan, Yi Gu, Victor Miller, Yonghao Zhuang, Guowei He, Haonan Li, Fajri Koto, Liping Tang, Nikhil Ranjan, Zhiqiang Shen, Xuguang Ren, Roberto Iriondo, Cun Mu, Zhiting Hu, Mark Schulze, Preslav Nakov, Tim Baldwin, and Eric P. Xing. LLM360: Towards fully transparent open-source LLMs. *CoRR*, abs/2312.06550, 2023. doi: 10.48550/ARXIV.2312.06550. URL <https://doi.org/10.48550/arXiv.2312.06550>.
- Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. In *7th International Conference on Learning Representations, ICLR 2019, New Orleans, LA, USA, May 6-9, 2019*. OpenReview.net, 2019. URL <https://openreview.net/forum?id=Bkg6RiCqY7>.
- Song Mei, Andrea Montanari, and Phan-Minh Nguyen. A mean field view of the landscape of two-layer neural networks. *Proceedings of the National Academy of Sciences*, 115(33):E7665–E7671, 2018. doi: 10.1073/pnas.1806579115. URL <https://www.pnas.org/doi/abs/10.1073/pnas.1806579115>.
- Naveen Mellempudi, Sudarshan Srinivasan, Dipankar Das, and Bharat Kaul. Mixed precision training with 8-bit floating point. *CoRR*, abs/1905.12334, 2019. URL <http://arxiv.org/abs/1905.12334>.
- Stephen Merity, Caiming Xiong, James Bradbury, and Richard Socher. Pointer sentinel mixture models. In *5th International Conference on Learning Representations, ICLR 2017, Toulon, France, April 24-26, 2017, Conference Track Proceedings*. OpenReview.net, 2017. URL <https://openreview.net/forum?id=Byj72udxe>.
- Paulius Micikevicius, Sharan Narang, Jonah Alben, Gregory F. Diamos, Erich Elsen, David García, Boris Ginsburg, Michael Houston, Oleksii Kuchaiev, Ganesh Venkatesh, and Hao Wu. Mixed precision training. In *6th International Conference on Learning Representations, ICLR 2018, Vancouver, BC, Canada, April 30 - May 3, 2018, Conference Track Proceedings*. OpenReview.net, 2018. URL <https://openreview.net/forum?id=rlgs9JgRZ>.
- Paulius Micikevicius, Dusan Stosic, Neil Burgess, Marius Cornea, Pradeep Dubey, Richard Grisenthwaite, Sangwon Ha, Alexander Heinecke, Patrick Judd, John Kamalu, Naveen Mellempudi, Stuart F. Oberman, Mohammad Shoeybi, Michael Y. Siu, and Hao Wu. FP8 formats for deep learning. *CoRR*, abs/2209.05433, 2022. doi: 10.48550/ARXIV.2209.05433. URL <https://doi.org/10.48550/arXiv.2209.05433>.
- Microsoft. Maximal update parametrization (μ P) and hyperparameter transfer (μ Transfer). <https://github.com/microsoft/mup>, 2024.
- Deepak Narayanan, Mohammad Shoeybi, Jared Casper, Patrick LeGresley, Mostofa Patwary, Vijay Korthikanti, Dmitri Vainbrand, Prethvi Kashinkunti, Julie Bernauer, Bryan Catanzaro, Amar Phanishayee, and Matei Zaharia. Efficient large-scale language model training on GPU clusters using Megatron-LM. In Bronis R. de Supinski, Mary W. Hall, and Todd Gamblin (eds.), *International Conference for High Performance Computing, Networking, Storage and Analysis, SC 2021, St. Louis, Missouri, USA, November 14-19, 2021*, pp. 58. ACM, 2021. doi: 10.1145/3458817.3476209. URL <https://doi.org/10.1145/3458817.3476209>.

- Lorenzo Noci, Sotiris Anagnostidis, Luca Biggio, Antonio Orvieto, Sidak Pal Singh, and Aurelien Lucchi. Signal propagation in transformers: Theoretical perspectives and the role of rank collapse. In S. Koyejo, S. Mohamed, A. Agarwal, D. Belgrave, K. Cho, and A. Oh (eds.), *Advances in Neural Information Processing Systems*, volume 35, pp. 27198–27211. Curran Associates, Inc., 2022. URL https://proceedings.neurips.cc/paper_files/paper/2022/file/ae0cba715b60c4052359b3d52a2cff7f-Paper-Conference.pdf.
- Badreddine Noune, Philip Jones, Daniel Justus, Dominic Masters, and Carlo Luschi. 8-bit numerical formats for deep neural networks. *CoRR*, abs/2206.02915, 2022. doi: 10.48550/ARXIV.2206.02915. URL <https://doi.org/10.48550/arXiv.2206.02915>.
- NVIDIA. Using FP8 with transformer engine. https://docs.nvidia.com/deeplearning/transformer-engine/user-guide/examples/fp8_primer.html, 2024a.
- NVIDIA. Transformer Engine. <https://github.com/NVIDIA/TransformerEngine>, 2024b.
- OpenAI. GPT-4 technical report. *CoRR*, abs/2303.08774, 2023. doi: 10.48550/ARXIV.2303.08774. URL <https://doi.org/10.48550/arXiv.2303.08774>.
- Houwen Peng, Kan Wu, Yixuan Wei, Guoshuai Zhao, Yuxiang Yang, Ze Liu, Yifan Xiong, Ziyue Yang, Bolin Ni, Jingcheng Hu, Ruihang Li, Miaosen Zhang, Chen Li, Jia Ning, Ruizhe Wang, Zheng Zhang, Shuguang Liu, Joe Chau, Han Hu, and Peng Cheng. FP8-LM: training FP8 large language models. *CoRR*, abs/2310.18313, 2023. doi: 10.48550/ARXIV.2310.18313. URL <https://doi.org/10.48550/arXiv.2310.18313>.
- Sergio P. Perez, Yan Zhang, James Briggs, Charlie Blake, Josh Levy-Kramer, Paul Balanca, Carlo Luschi, Stephen Barlow, and Andrew Fitzgibbon. Training and inference of large language models using 8-bit floating point. *CoRR*, abs/2309.17224, 2023. doi: 10.48550/ARXIV.2309.17224. URL <https://doi.org/10.48550/arXiv.2309.17224>.
- Ben Poole, Subhaneil Lahiri, Maithra Raghu, Jascha Sohl-Dickstein, and Surya Ganguli. Exponential expressivity in deep neural networks through transient chaos. In D. Lee, M. Sugiyama, U. Luxburg, I. Guyon, and R. Garnett (eds.), *Advances in Neural Information Processing Systems*, volume 29. Curran Associates, Inc., 2016. URL https://proceedings.neurips.cc/paper_files/paper/2016/file/148510031349642de5ca0c544f31b2ef-Paper.pdf.
- Shikai Qiu, Andres Potapczynski, Marc Finzi, Micah Goldblum, and Andrew Gordon Wilson. Compute better spent: Replacing dense layers with structured matrices. *CoRR*, abs/2406.06248, 2024. doi: 10.48550/ARXIV.2406.06248. URL <https://doi.org/10.48550/arXiv.2406.06248>.
- Samyam Rajbhandari, Jeff Rasley, Olatunji Ruwase, and Yuxiong He. ZeRO: Memory optimizations toward training trillion parameter models. In Christine Cuicchi, Irene Qualters, and William T. Kramer (eds.), *Proceedings of the International Conference for High Performance Computing, Networking, Storage and Analysis, SC 2020, Virtual Event / Atlanta, Georgia, USA, November 9-19, 2020*, pp. 20. IEEE/ACM, 2020. doi: 10.1109/SC41405.2020.00024. URL <https://doi.org/10.1109/SC41405.2020.00024>.
- Prajit Ramachandran, Barret Zoph, and Quoc V. Le. Searching for activation functions. In *6th International Conference on Learning Representations, ICLR 2018, Vancouver, BC, Canada, April 30 - May 3, 2018, Workshop Track Proceedings*. OpenReview.net, 2018. URL <https://openreview.net/forum?id=Hkuq2EkPf>.
- Samuel S. Schoenholz, Justin Gilmer, Surya Ganguli, and Jascha Sohl-Dickstein. Deep information propagation. In *5th International Conference on Learning Representations, ICLR 2017, Toulon, France, April 24-26, 2017, Conference Track Proceedings*. OpenReview.net, 2017. URL <https://openreview.net/forum?id=H1W1UN9gg>.
- Noam Shazeer. GLU variants improve transformer. *CoRR*, abs/2002.05202, 2020. URL <https://arxiv.org/abs/2002.05202>.

- Zhiqiang Shen, Tianhua Tao, Liqun Ma, Willie Neiswanger, Zhengzhong Liu, Hongyi Wang, Bowen Tan, Joel Hestness, Natalia Vassilieva, Daria Soboleva, and Eric P. Xing. SlimPajama-DC: Understanding data combinations for LLM training. *CoRR*, abs/2309.10818, 2023. doi: 10.48550/ARXIV.2309.10818. URL <https://doi.org/10.48550/arXiv.2309.10818>.
- Jianlin Su, Murtadha H. M. Ahmed, Yu Lu, Shengfeng Pan, Wen Bo, and Yunfeng Liu. Roformer: Enhanced transformer with rotary position embedding. *Neurocomputing*, 568:127063, 2024. doi: 10.1016/J.NEUCOM.2023.127063. URL <https://doi.org/10.1016/j.neucom.2023.127063>.
- Mingjie Sun, Xinlei Chen, J. Zico Kolter, and Zhuang Liu. Massive activations in large language models. *CoRR*, abs/2402.17762, 2024. doi: 10.48550/ARXIV.2402.17762. URL <https://doi.org/10.48550/arXiv.2402.17762>.
- Xiao Sun, Jungwook Choi, Chia-Yu Chen, Naigang Wang, Swagath Venkataramani, Vijayalakshmi Srinivasan, Xiaodong Cui, Wei Zhang, and Kailash Gopalakrishnan. Hybrid 8-bit floating point (HFP8) training and inference for deep neural networks. In Hanna M. Wallach, Hugo Larochelle, Alina Beygelzimer, Florence d’Alché-Buc, Emily B. Fox, and Roman Garnett (eds.), *Advances in Neural Information Processing Systems 32: Annual Conference on Neural Information Processing Systems 2019, NeurIPS 2019, December 8-14, 2019, Vancouver, BC, Canada*, pp. 4901–4910, 2019. URL <https://proceedings.neurips.cc/paper/2019/hash/65fc9fb4897a89789352e211ca2d398f-Abstract.html>.
- Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurélien Rodriguez, Armand Joulin, Edouard Grave, and Guillaume Lample. Llama: Open and efficient foundation language models. *CoRR*, abs/2302.13971, 2023a. doi: 10.48550/ARXIV.2302.13971. URL <https://doi.org/10.48550/arXiv.2302.13971>.
- Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajwal Bhargava, Shruti Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton-Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony Hartshorn, Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa, Isabel Kloumann, Artem Korenev, Punit Singh Koura, Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov, Pushkar Mishra, Igor Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashi Rungta, Kalyan Saladi, Alan Schelten, Ruan Silva, Eric Michael Smith, Ranjan Subramanian, Xiaoqing Ellen Tan, Binh Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu, Zheng Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan, Melanie Kambadur, Sharan Narang, Aurélien Rodriguez, Robert Stojnic, Sergey Edunov, and Thomas Scialom. Llama 2: Open foundation and fine-tuned chat models. *CoRR*, abs/2307.09288, 2023b. doi: 10.48550/ARXIV.2307.09288. URL <https://doi.org/10.48550/arXiv.2307.09288>.
- Naigang Wang, Jungwook Choi, Daniel Brand, Chia-Yu Chen, and Kailash Gopalakrishnan. Training deep neural networks with 8-bit floating point numbers. In Samy Bengio, Hanna M. Wallach, Hugo Larochelle, Kristen Grauman, Nicolò Cesa-Bianchi, and Roman Garnett (eds.), *Advances in Neural Information Processing Systems 31: Annual Conference on Neural Information Processing Systems 2018, NeurIPS 2018, December 3-8, 2018, Montréal, Canada*, pp. 7686–7695, 2018. URL <https://proceedings.neurips.cc/paper/2018/hash/335d3d1cd7ef05ec77714a215134914c-Abstract.html>.
- Xi Wang and Laurence Aitchison. How to set adamw’s weight decay as you scale model and dataset size. *CoRR*, abs/2405.13698, 2024. doi: 10.48550/ARXIV.2405.13698. URL <https://doi.org/10.48550/arXiv.2405.13698>.
- Mitchell Wortsman, Peter J. Liu, Lechao Xiao, Katie Everett, Alex Alemi, Ben Adlam, John D. Co-Reyes, Izzeddin Gur, Abhishek Kumar, Roman Novak, Jeffrey Pennington, Jascha Sohl-Dickstein, Kelvin Xu, Jaehoon Lee, Justin Gilmer, and Simon Kornblith. Small-scale proxies for large-scale transformer training instabilities. *CoRR*, abs/2309.14322, 2023. doi: 10.48550/ARXIV.2309.14322. URL <https://doi.org/10.48550/arXiv.2309.14322>.

- xAI. Grok-1. <https://github.com/xai-org/grok-1>, 2024.
- 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, Jianxin Yang, Jin Xu, Jingren Zhou, Jinze Bai, Jinzheng He, Junyang Lin, Kai Dang, Keming Lu, Keqin Chen, Kexin Yang, Mei Li, Mingfeng Xue, Na Ni, Pei Zhang, Peng Wang, Ru Peng, Rui Men, Ruize Gao, Runji Lin, Shijie Wang, Shuai Bai, Sinan Tan, Tianhang Zhu, Tianhao Li, Tianyu Liu, Wenbin Ge, Xiaodong Deng, Xiaohuan Zhou, Xingzhang Ren, Xinyu Zhang, Xipin Wei, Xuancheng Ren, Xuejing Liu, Yang Fan, Yang Yao, Yichang Zhang, Yu Wan, Yunfei Chu, Yuqiong Liu, Zeyu Cui, Zhenru Zhang, Zhifang Guo, and Zhihao Fan. Qwen2 technical report. *CoRR*, abs/2407.10671, 2024. doi: 10.48550/ARXIV.2407.10671. URL <https://doi.org/10.48550/arXiv.2407.10671>.
- Greg Yang. Tensor programs I: Wide feedforward or recurrent neural networks of any architecture are Gaussian processes. *CoRR*, abs/1910.12478, 2019. URL <http://arxiv.org/abs/1910.12478>.
- Greg Yang. Tensor programs II: Neural tangent kernel for any architecture. *CoRR*, abs/2006.14548, 2020a. URL <https://arxiv.org/abs/2006.14548>.
- Greg Yang. Tensor programs III: Neural matrix laws. *CoRR*, abs/2009.10685, 2020b. URL <https://arxiv.org/abs/2009.10685>.
- Greg Yang and Edward J. Hu. Tensor programs IV: Feature learning in infinite-width neural networks. In Marina Meila and Tong Zhang (eds.), *Proceedings of the 38th International Conference on Machine Learning, ICML 2021, 18-24 July 2021, Virtual Event*, volume 139 of *Proceedings of Machine Learning Research*, pp. 11727–11737. PMLR, 2021. URL <http://proceedings.mlr.press/v139/yang21c.html>.
- Greg Yang and Etai Littwin. Tensor programs IIb: Architectural universality of neural tangent kernel training dynamics. In Marina Meila and Tong Zhang (eds.), *Proceedings of the 38th International Conference on Machine Learning, ICML 2021, 18-24 July 2021, Virtual Event*, volume 139 of *Proceedings of Machine Learning Research*, pp. 11762–11772. PMLR, 2021. URL <http://proceedings.mlr.press/v139/yang21f.html>.
- Greg Yang and Etai Littwin. Tensor programs IVb: Adaptive optimization in the infinite-width limit. *CoRR*, abs/2308.01814, 2023. doi: 10.48550/ARXIV.2308.01814. URL <https://doi.org/10.48550/arXiv.2308.01814>.
- 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. *CoRR*, abs/2203.03466, 2022. doi: 10.48550/ARXIV.2203.03466. URL <https://doi.org/10.48550/arXiv.2203.03466>.
- Greg Yang, James B. Simon, and Jeremy Bernstein. A spectral condition for feature learning. *CoRR*, abs/2310.17813, 2023a. doi: 10.48550/ARXIV.2310.17813. URL <https://doi.org/10.48550/arXiv.2310.17813>.
- Greg Yang, Dingli Yu, Chen Zhu, and Soufiane Hayou. Tensor programs VI: Feature learning in infinite-depth neural networks. *CoRR*, abs/2310.02244, 2023b. doi: 10.48550/ARXIV.2310.02244. URL <https://doi.org/10.48550/arXiv.2310.02244>.
- Chao Yu and Zhiguo Su. Symmetrical Gaussian error linear units (SGELUs). *CoRR*, abs/1911.03925, 2019. URL <http://arxiv.org/abs/1911.03925>.
- Biao Zhang and Rico Sennrich. Root mean square layer normalization. In Hanna M. Wallach, Hugo Larochelle, Alina Beygelzimer, Florence d’Alché-Buc, Emily B. Fox, and Roman Garnett (eds.), *Advances in Neural Information Processing Systems 32: Annual Conference on Neural Information Processing Systems 2019, NeurIPS 2019, December 8-14, 2019, Vancouver, BC, Canada*, pp. 12360–12371, 2019. URL <https://proceedings.neurips.cc/paper/2019/hash/1e8a19426224ca89e83cef47f1e7f53b-Abstract.html>.

A ADDITIONAL EXPERIMENTAL DETAILS

A.1 EXPERIMENTAL SETUP

Our experimental analysis of u- μ P was conducted by adapting the codebase used for Tensor Programs V, allowing us to compare μ P and u- μ P in the same setting. We change various experimental settings from the μ P paper to make our experiments better reflect standard training procedures, particularly the dataset which we switch from WikiText-2 to the larger WikiText-103 (Merity et al., 2017). Where not specified otherwise, the default setting used in our experiments are given in Table 5. These also represent the settings of our proxy model.

Dataset	WikiText-103 (Merity et al., 2017)
Sequence length	256
Vocab size	32000
Training set tokens	138M
Architecture	Llama (Touvron et al., 2023a) (Transformer, PreNorm, RMSNorm, SwiGLU, RoPE, “untied” embeddings), non-trainable RMSNorm parameters.
Width	256 (scaled up to 4096)
Depth	4
Number of heads	4 (scaled up to 64)
Head dimension	64
Total parameters	19.5M (scaled up to 1.07B)
Batch size	64
Training steps	8192 (0.97 epochs)
LR schedule	Cosine to 10%, 2000 steps warm-up
Optimizer	AdamW ($\beta_1, \beta_2, \epsilon$) = (0.9, 0.999, 10^{-8})
Weight decay	2^{-13} , independent (Loshchilov & Hutter, 2019)
Dropout	0.0
μ P HP search range	$\eta \in [2^{-10}, 2^{-6}]$ $\hat{\eta}_{\text{emb}} \in [2^0, 2^8]$
u- μ P HP search range	$\sigma_{\text{init}}, \alpha_{\text{emb}}, \alpha_{\text{attn}}, \alpha_{\text{output}} \in [2^{-2}, 2^2]$ $\eta \in [2^{-1}, 2^3]$ $\alpha_{\text{attn}} \in [2^{-2}, 2^2]$ $\alpha_{\text{residual}}, \alpha_{\text{residual-attn-ratio}}, \alpha_{\text{ffn-act}}, \alpha_{\text{output}} \in [2^{-3}, 2^3]$
μ P HP defaults	$\sigma_{\text{init}} = \alpha_{\text{emb}} = \alpha_{\text{attn}} = \alpha_{\text{output}} = \hat{\eta}_{\text{emb}} = 1$
u- μ P HP defaults	$\alpha_{\text{residual}} = \alpha_{\text{residual-attn-ratio}} = \alpha_{\text{ffn-act}} = \alpha_{\text{output}} = \alpha_{\text{attn}} = 1$

Table 5: Default hyperparameters and training settings.

A.2 VALIDATING OUR EXPERIMENTAL SETUP

In this section we run a series of ablations to validate decisions relating to our experimental setup given above. In particular, we examine the effect of using repeated data, the effect of using a shorter warmup duration, and the effect of different final learning rates at the end of decay.

A.2.1 REPEATED DATA

As outlined in Table 5, our standard training setup uses 0.97 epochs of the WikiText-103 dataset (50x larger than the WikiText-2 dataset used in Tensor Programs V). However on our batch size and training steps scaling experiments in Figure 4 we train on up to $4\times$ the amount of data than in our standard setup, and hence use up to 4 epochs.

Though this is still a small level of repeated data, this moves our training slightly into the over-fitting regime. Based on this change, we here investigate the hypothesis that this regime has better or worse transfer of the optimal LR than the non-overfitting regime, and hence our results could be misleading. To do so, we repeated these experiments with the same number of tokens, but using the much larger SlimPajama dataset (Shen et al., 2023) where we use < 1 epoch.

The results for this experiment are seen in Figure 7. The shape of curves is very similar across the two datasets, for both batch size and training steps (albeit with a higher loss, due to the more varied nature of SlimPajama). From this we conclude that the effect of repeated data from our use of WikiText-103 is not significant.

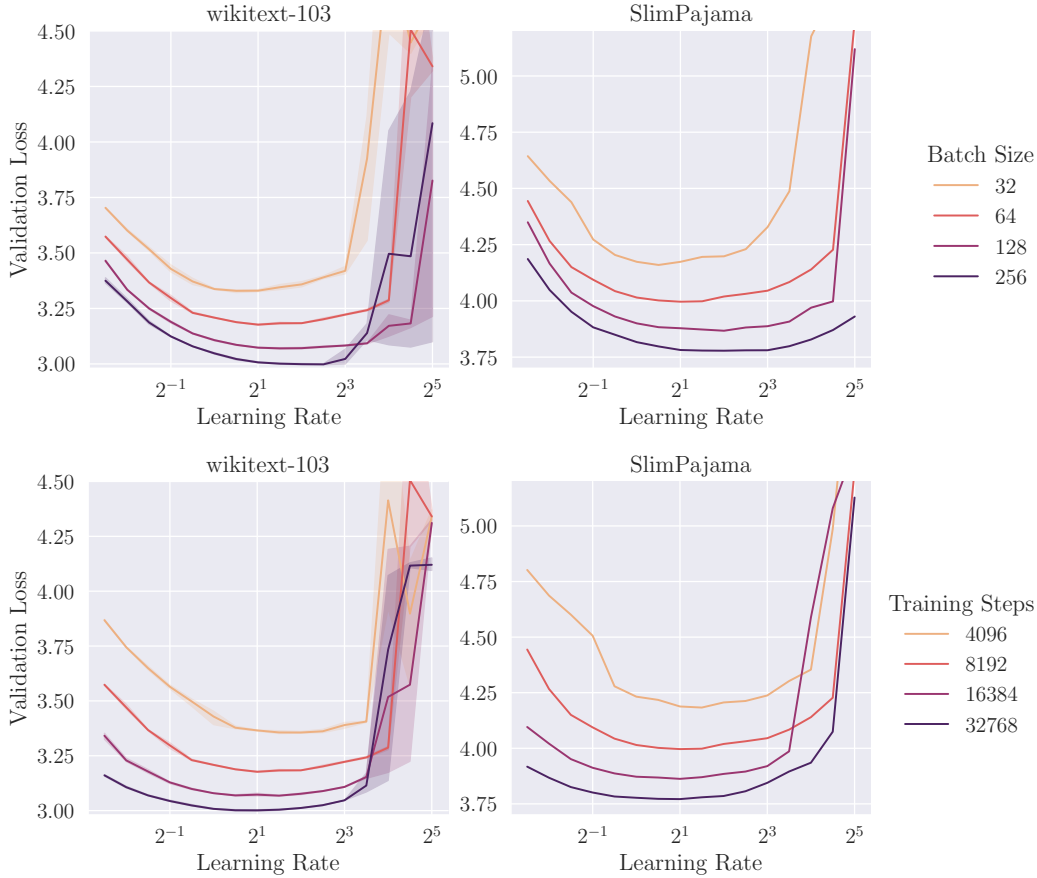


Figure 7: A repeat of the batch size and training steps experiments in Figure 4, but using the larger SlimPajama dataset where no data is repeated. In both settings our validation loss basins take the same shape, indicating that our analysis using the WikiText-103 dataset holds.

A.2.2 WARMUP DURATION

For our experimental setup (Table 5) we use a longer duration of warmup than in our large-scale setup (Table 6). We do so out of caution, as we use fewer tokens-per-batch for the smaller-scale experiments and so may require longer warmup. However, doing so also creates the risk of spending too large a proportion of training doing warmup, which could affect transfer.

To investigate this effect, we run two experiments. Firstly, we re-run the experiment for LR transfer over training steps, shown in Figure 4, on a quarter of the warmup steps. This is shown in Figure 8 (left). The main effect appears to be higher loss for larger learning rates, but the optima are largely unchanged. The only exception is the 4096-step run, where the optimum shifts left and the loss improves slightly. This appears to now align the optimum better with the other training durations, but leads to narrower basins as a result, suggesting a trade-off for this particular experiment.

However, all our other experimental runs use the 8192-step configuration, which has a consistent optimum regardless of warmup duration here. To investigate the effect of reduced warmup on width transfer at this particular step-count, we re-run our experiment in Figure 1 (b) under the shorter warmup duration, shown in Figure 8 (right). The only significant impact of this change is to narrow the basins, inducing no significant change in the optimal LR. As such, we conclude that using 2000 steps of warmup in our experimental setup is a reasonable choice, and both give the same width transfer.

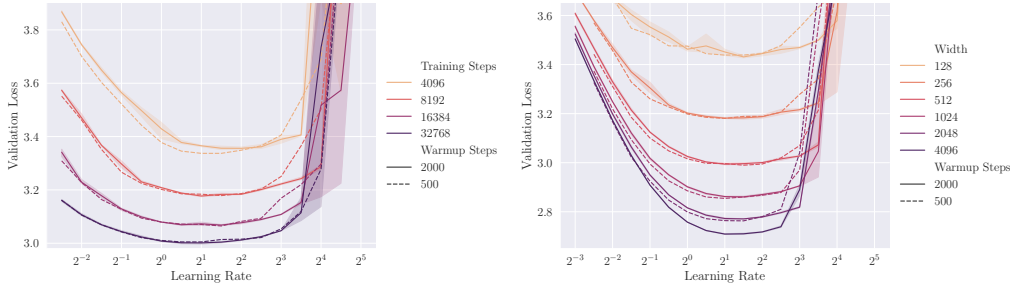


Figure 8: (Left) Learning rate transfer across training steps under different numbers of warmup steps. (Right) Learning rate transfer across width under different numbers of warmup steps. In this setting (training steps = 8192) the optimal LR is consistent, meaning either warmup regime can be used, though the longer gives wider basins.

A.2.3 LEARNING RATE DECAY TARGET

In all our experiments we use a cosine decay of our learning rate down to 10% of the maximum. This follows the standard approach taken by most LLM training projects (Touvron et al., 2023b; Biderman et al., 2023; Groeneveld et al., 2024; Almazrouei et al., 2023; Yang et al., 2024). However, recent research has indicated that this may not be the optimal decay target, with implications for LR transfer. Hägele et al. (2024) show that the choice of target percentage can alter the shape of transfer curves and potentially shift the optimum value (Figure 21, right). They also suggest that using a fixed target value may work better than a percentage (Figure 22, right), which could be swept separately. Anonymous (2024) separately suggest that linear decay to zero is the most effective scheme.

Though using the optimal decay scheme is not necessarily essential to the validity of our method, any implications of different schemes on transfer properties should be investigated. To do so, we run two experiments. The first sweeps the learning rate for our standard model at various percentages and fixed values of cosine decay target, including zero, in Figure 9 (left). Lower decay targets perform better here, including zero, suggesting that this simple rule may be ideal.

We then re-run our width transfer experiment from Figure 1 (b) but with our LR decaying to 0, and plot the result in Figure 9 (right). This leads to slightly better results for large learning rates, though for large models this difference diminishes. Fortunately the effect this decay target has on the shape of curves (and hence optimal LR transfer) is minimal, indicating that our conclusions are not effected significantly by the choice of decay target.

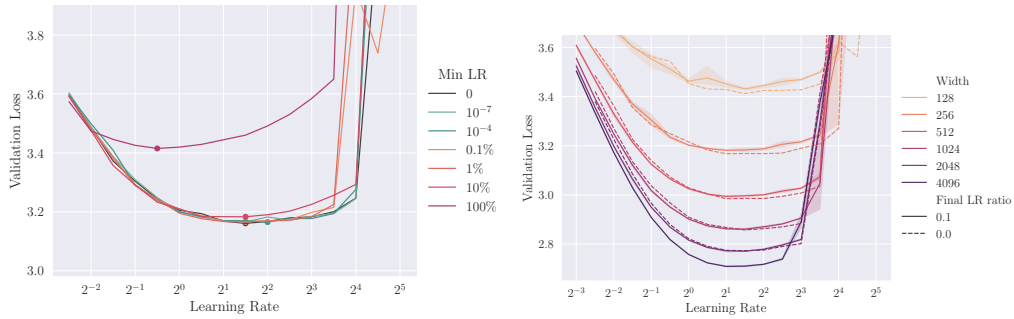


Figure 9: (Left) A learning rate sweep over LR targets of different types (percentage, fixed and zero) on our standard model. (Right) Using the zero and 10% learning rate targets, LR transfer over width.

A.3 PER-TENSOR LEARNING RATES

In Section 3.3 we relax the requirement for each weight tensor in the u- μ P model to have an associated tuneable learning-rate multiplier on top of the global learning rate. Whilst this does reduce the theoretical expressivity of the u- μ P scheme, Figure 10 shows that using a single globally optimized learning rate is already at or close to the optimal choice for all weight tensors, and therefore it is reasonable to drop these multipliers in favor of reducing the number of HPs. However, a practitioner attempting to absolutely maximize the task performance of their model could experiment with tuning a few key per-tensor LRs, in particular the embedding table.



Figure 10: Independently varying per-tensor learning rate multipliers η_W , using the u- μ P model of width 256 from Figure 1 with optimized global learning rate $2^{1.5}$ as the starting point. Where applicable, the same multiplier is used for tensors of the same name across transformer layers. Each subplot fixes all but one multiplier at 1, therefore the midpoint of each subplot is precisely the u- μ P₂₅₆ model from Figure 1.

A.4 HYPERPARAMETER INDEPENDENCE

In Section 4.1 we explore the question of HP independence under μ P and u- μ P. The following plots in Figures 11 and 12 show the result of a 2D sweep over every pair of HPs under each scheme. All other HPs are held at 1 when not swept, except the η which is held at $2^{-7.5}$ for μ P and $2^{1.5}$ for u- μ P, and $\hat{\eta}_{\text{emb}}$ which is held at 2^4 for μ P.

These results show visual dependence between μ P hyperparameters as a diagonal structure in the grids, such as $(\hat{\eta}_{\text{emb}}, \sigma_{\text{init}})$ and $(\eta, \alpha_{\text{attn}})$. We quantify this in the plot in Figure 3, where we use a measure of HP dependence termed transfer error. This is explained verbally in Section 4.1, and we provide an algorithmic description in Algorithm 1. We note that differences in transfer error between the two methods may also be influenced by the flatness of the optimum. The HP and loss values used for our transfer error calculations are those in Figures 11 and 12.

Algorithm 1 Transfer Error

Require: A ‘fixed’ HP with candidate values $F = \{f_1, \dots, f_n\}$, a ‘transfer’ HP with candidate values $T = \{t_1, \dots, t_m\}$, a function that gives the final validation loss for the pair of HPs $L : F \times T \rightarrow \mathbb{R}$ (assuming all other HPs are fixed at default values).

```

err  $\leftarrow$  0
 $f^*, t^* \leftarrow \text{argmin}(L)$ 
for  $f$  in  $F$  do
  if  $f \neq f^*$  then
     $t \leftarrow \text{argmin}(L(f))$ 
    err  $+= L(f^*, t) - L(f^*, t^*)$ 
  end if
end for
return err / ( $n - 1$ )

```

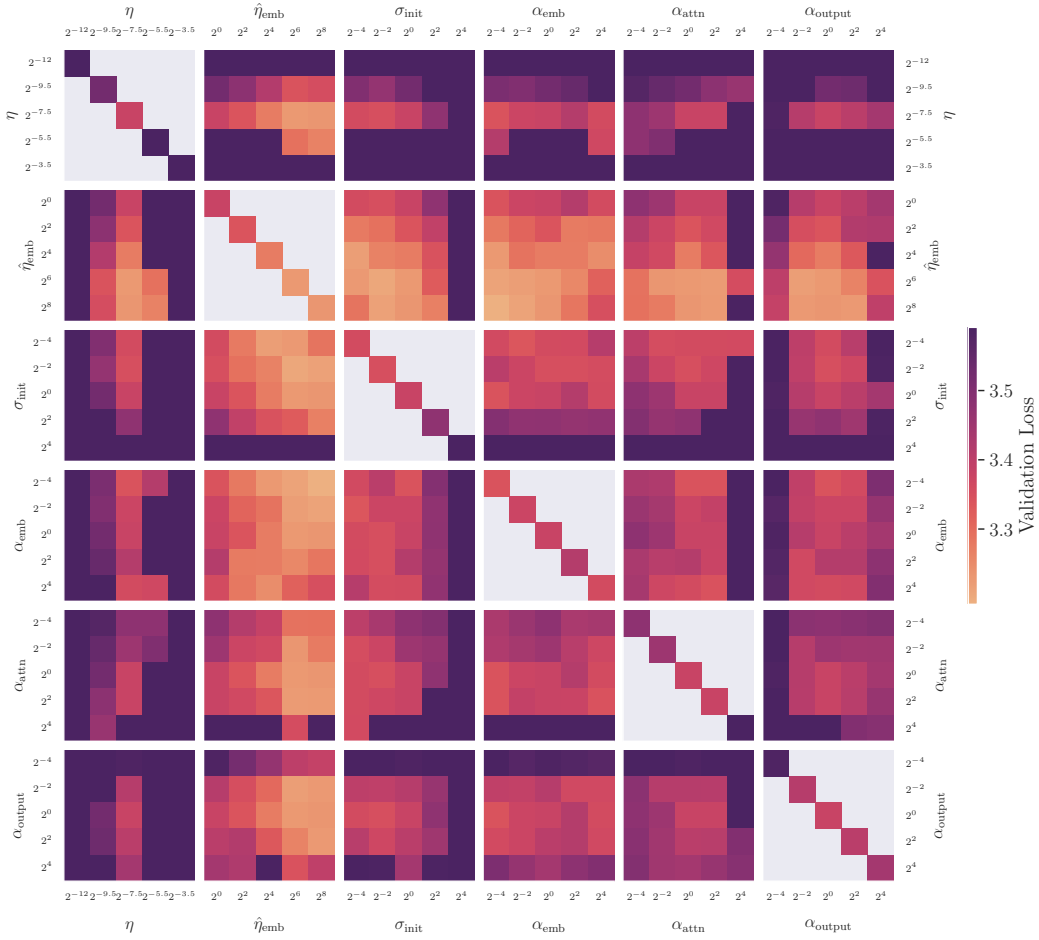


Figure 11: Hyperparameter coupling sweep for μP . Note strong coupling between optima, e.g. in the cases of $(\hat{\eta}_{\text{emb}}, \sigma_{\text{init}})$ and $(\eta, \alpha_{\text{attn}})$. See also: u- μP , Figure 12.

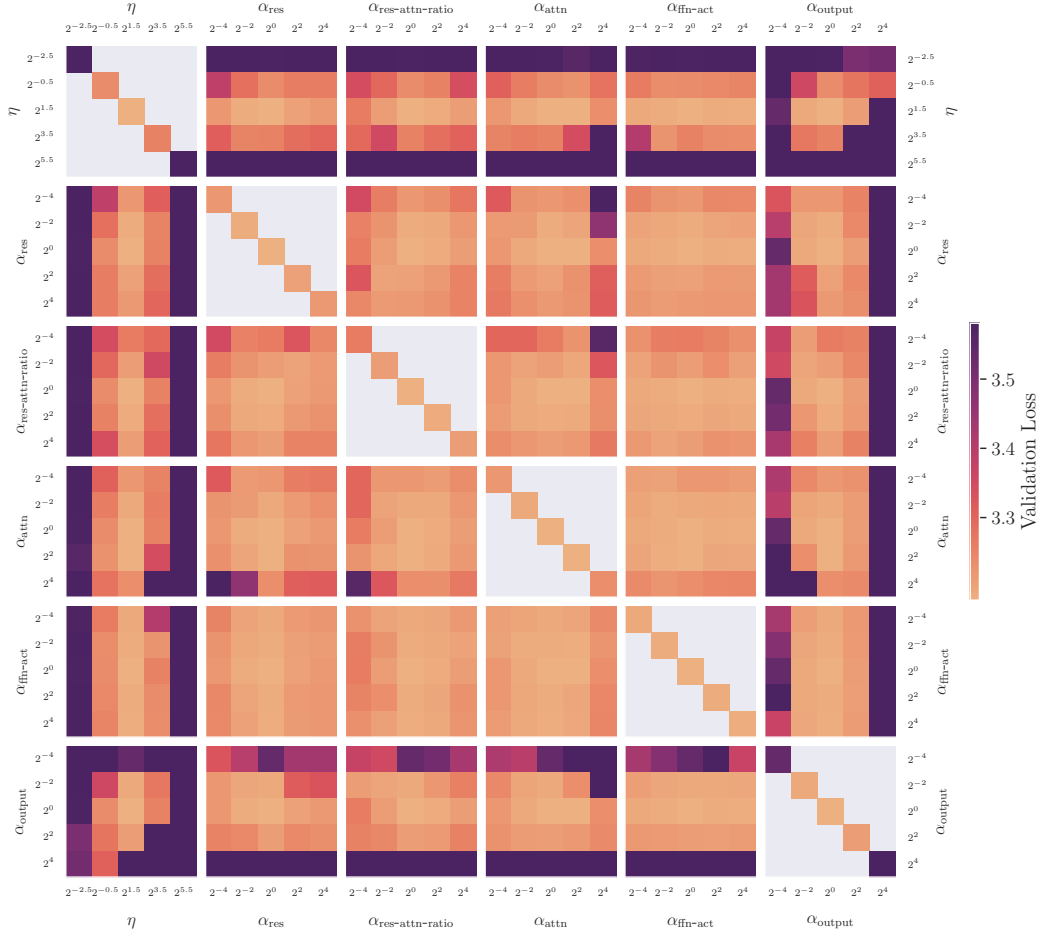


Figure 12: Hyperparameter coupling sweep for u- μ P. Note less coupling than with μ P, see Figure 11.

A.5 HYPERPARAMETER SEARCH

Here we outline the particular search processes used for our μP and $\text{u-}\mu\text{P}$ HP sweeps in Figure 1 (a). The *random search* samples uniformly from a grid defined over all *extended* HPs (extended HP sets are defined in Table 3, with grid values defined in Table 5). We perform the random search over 339 runs, each of which is a full training of the width-256 proxy model. We then simulate the effect of shorter searches at various run-counts by taking a random sample of the results, resulting in the smooth curve over run-count shown.

The *independent search* consists of the following phases:

1. Perform a 1D line search for an optimal learning rate, with other hyperparameters set to their default values (9 runs).
2. For each hyperparameter in parallel, perform a 1D line search (330 runs).
3. Combine the best settings from step 2, and re-evaluate (6 runs).

The number of runs in the 1D line search is an order of magnitude higher than is required in practice. We do so to form a fair comparison with the random search, which benefits from this large number of runs. The number of runs for the 1D line search could be reduced further by using binary search, though this would require sequential runs and limit the extent of parallelism.

A.6 HYPERPARAMETER TRANSFER EXPERIMENTS

Baseline μP transfer Figure 13 is a companion plot to Figure 4 in the body of the paper, showing the LR transfer of the baseline μP model over the same axes. $\text{u-}\mu\text{P}$ shows marginally more stable HP transfer here relative to the baseline, and at a consistently lower loss.

LR transfer over width The transfer experiments shown in Figure 1 (b) use the non-LR HPs found in Figure 1 (a) (indicated by the circled points), rather than using default HP values. For the $\text{u-}\mu\text{P}$ sweep we take the HPs at the end of the LR portion of the independent search, as these are already close-to-optimal, and means only 9 runs were required in the sweep. In contrast, for μP it is necessary to use the results of the random search over a large number of runs.

LR transfer over other axes For the training steps, batch size and depth transfer experiments in Figure 4, all HP values are fixed to 1 except LR which is swept. As with width transfer, $\text{u-}\mu\text{P}$ outperforms μP here using these default HP values. Reducing training steps is done by fixing the number of warm-up steps (at 2000) and still cosine-decaying the learning rate to 10%; all that changes is the number of post-warm-up steps. We found this to be more effective than cutting-short the decay schedule. For both Figure 1 (b) and Figure 4 we sweep the LR over a logarithmically-spaced grid of step $2^{1/2} \times$, with 3 runs for each point.

Additionally, in Figure 14 we show learning rate transfer over sequence length for both μP and $\text{u-}\mu\text{P}$ fixing either tokens per batch or sequences per batch. In both scenarios $\text{u-}\mu\text{P}$ shows not only better absolute training performance, but also better transfer behaviour as sequence length increases. Since our default proxy sequence length is 256, using μP to transfer to sequence length 2048 would result in minimal improvements or even a degradation in validation loss, whereas the $\text{u-}\mu\text{P}$ shows much greater and more consistent improvements.

Other HP transfer over width For our non-LR HP transfer results in Figure 15, we note that good transfer under μP has not been demonstrated for all HPs used in the literature. This is particularly true for the $\hat{\eta}_{\text{emb}}$ HP, which has poor transfer under μP . Our investigation here led to our identification of the need to adjust the embedding LR scaling rule outlined in Section 3.4. In many cases users have not swept this HP, but instead swept the corresponding parameter multiplier α_{emb} . How this HP interacts with the embedding LR scaling problem identified (and our proposed fix) remains to be explored, though we note in Figure 15 that it also appears to have poor transfer.

Combined HP transfer Whilst Figure 15 demonstrates the transfer of individual hyperparameters over width, Figure 16 instead demonstrates the simultaneous transfer of all hyperparameters when co-optimized on the small-scale proxy model, as is done for $\mu\text{Transfer}$. The μP and $\text{u-}\mu\text{P}$ points are

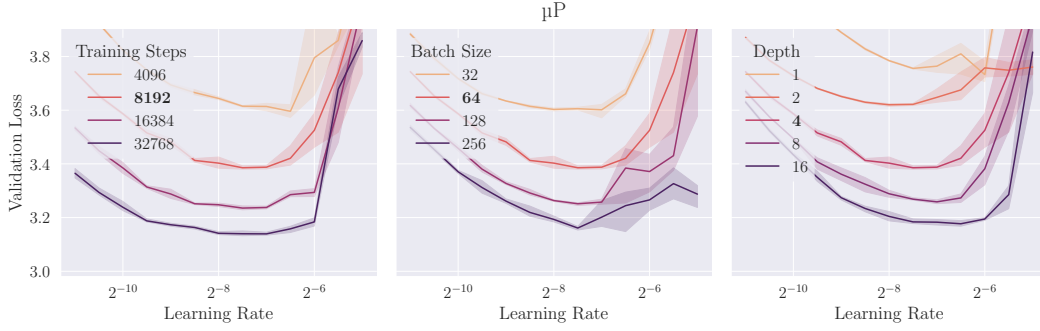


Figure 13: Learning rate transfer for μP over training steps, batch size and depth. This is a companion to the equivalent $u\text{-}\mu P$ plot in Figure 4.

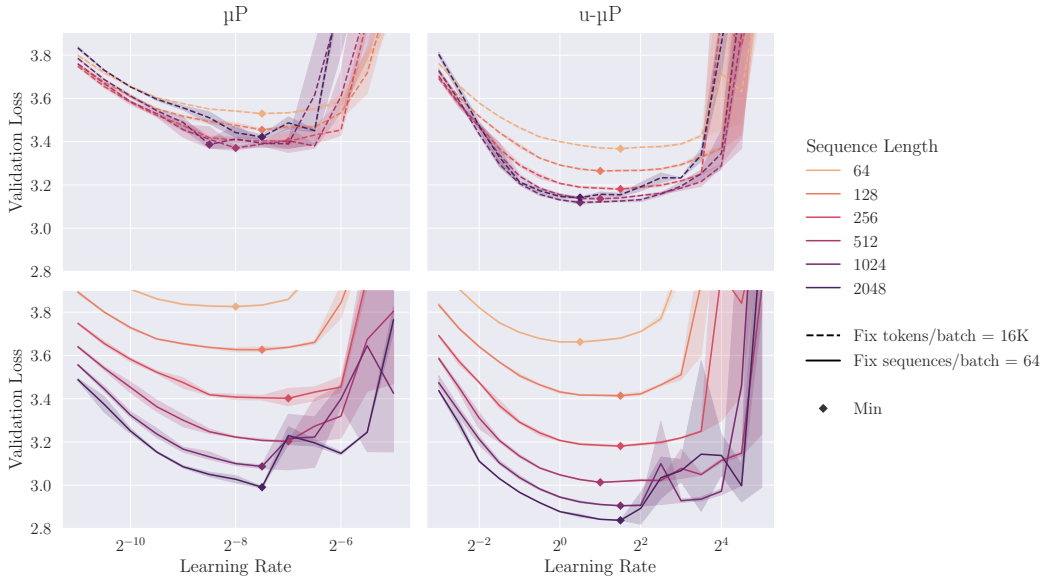


Figure 14: Transfer of learning rate over sequence length for μP (left) and $u\text{-}\mu P$ (right). As sequence length varies, we can fix the number of tokens per batch by inversely varying the number of sequences per batch (top). Alternatively we can fix the sequences per batch and allow the number of tokens per batch to vary with sequence length (bottom). In the latter case, larger sequence lengths mean the model sees more tokens during training, though as per Table 5 this translates to >1 epoch on WikiText-103 when sequence length goes above 256.

taken from Figure 1, with hyperparameters swept on a model of width 256 using a full random HP search and a simple learning rate sweep for μP and $u\text{-}\mu P$ respectively. The Standard Parametrization scheme, as shown in Table 3 requires choosing a learning rate and a weight-initialization scheme. We follow the initialization scheme of Pythia (Biderman et al., 2023), and transfer learning rate using a heuristic scaling factor of $\text{base-width}/\text{width}$, as is done in Dubey et al. (2024).

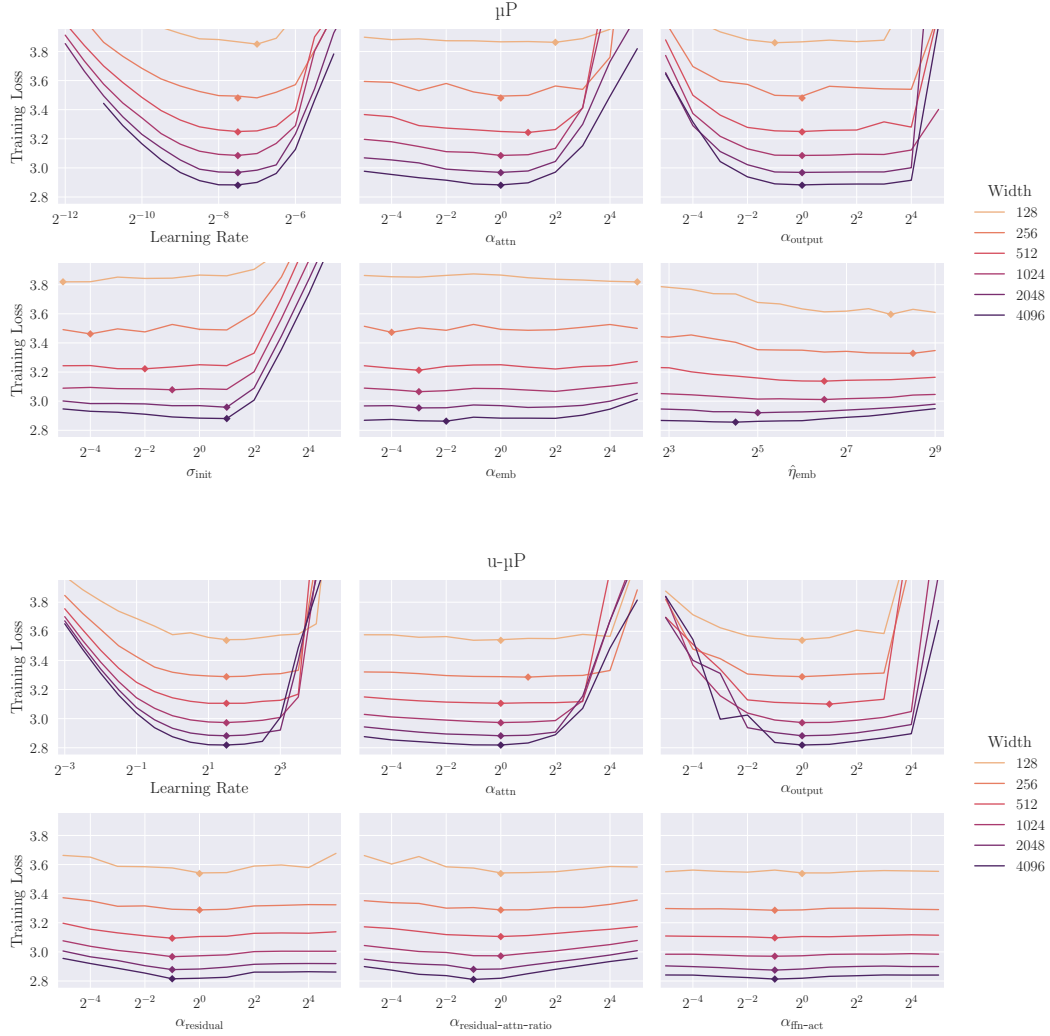


Figure 15: Transfer of model hyperparameters over width for μP (top) and $u\text{-}\mu\text{P}$ (bottom). When one hyperparameter is being swept, all others are fixed at 1, with the exception of Learning Rate $\eta = (2^{1.5}, 2^{-7.5})$ for ($u\text{-}\mu\text{P}$, μP).

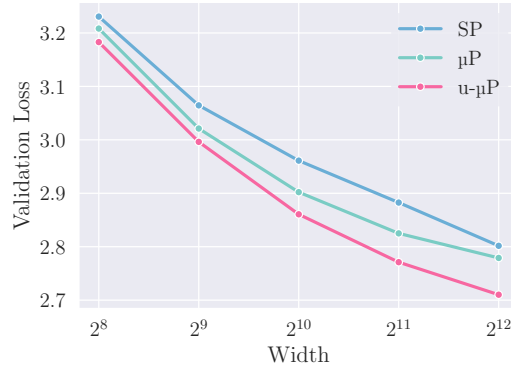


Figure 16: Transferring hyperparameters from width 256 up to 4096 using three different hyperparametrization schemes. μP and $u\text{-}\mu\text{P}$ results are as seen in Figure 1, whilst Standard Parametrization follows the initialization approach of Pythia (Biderman et al., 2023).

A.7 NUMERICAL PROPERTIES

Our analysis of the numerical properties of u- μ P focuses on the RMS of tensors that we wish to cast to FP8: linear module input activations, weights and output gradients. From the RMS training statistics plots in Figure 5 and Figure 17 we note that

1. μ P has gradients and weights with low RMS, at risk of FP8 underflow, whereas u- μ P starts with $\text{RMS} \approx 1$.
2. Many input activations do not grow RMS during training (due to a preceding non-trainable RMSNorm), however the attention out projection and FFN down projection have unconstrained input activations that grow considerably during training.
3. The decoder weight grows during training. Since it is preceded by a RMSNorm, the model may require scale growth in order to increase the scale of softmax inputs. Other weights grow slightly during training.
4. Gradients grow quickly but stabilize, except for attention out projection and FFN down projection, whose gradients shrink as the inputs grow.

We also evaluate how RMS growth is affected by model and training hyperparameters in the tensors that showed the highest end-training RMS, shown in Figure 18. This shows that the main parameter affecting scale growth is learning rate, with end-training RMS increasing to the right of the optimal LR basin, as training becomes unstable. End-training RMS is remarkably stable as width, depth, training steps and batch size are independently increased.

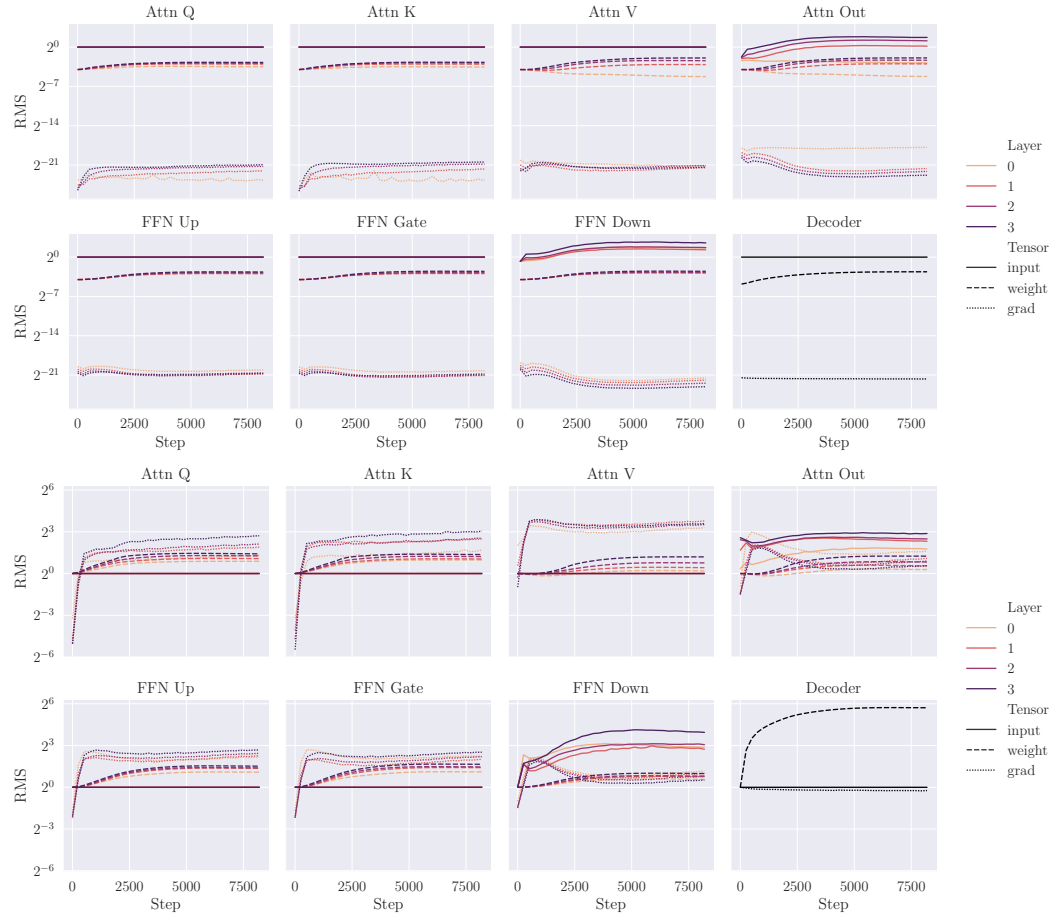


Figure 17: RMS during training, for all parametrized matmul inputs, for μ P (top) and u- μ P (bottom). Model width 256, default hyperparameters, $\eta = (2^1, 2^{-8})$ for (u- μ P, μ P).

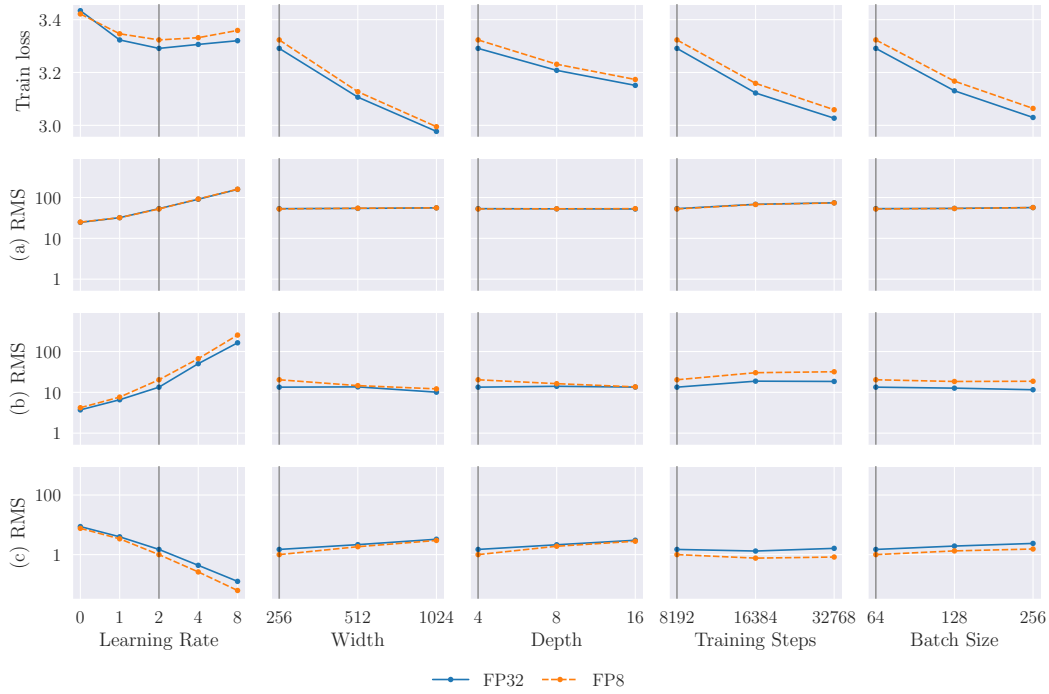


Figure 18: The effect of hyperparameters on FP8 training loss and on the end-training RMS of critical tensors: (a) decoder weight, (b) last-layer FFN down-projection input and (c) last-layer FFN down-projection output gradient. Only learning rate has a substantial effect on the end-training RMS. Vertical lines show the default setting of that hyperparameter, as used for all other plots.

A.8 LARGE-SCALE TRAINING DETAILS

Our large-scale training settings are given in Table 6. These are largely the same as our standard experiments (Table 5), but with many more tokens used for training and scaling up to a larger model-size.

Dataset	SlimPajama (Shen et al., 2023)
Sequence length	4096
Vocab size	65536
Training set tokens	600B
Architecture	Llama (Touvron et al., 2023a) (Transformer, PreNorm, RMSNorm, SwiGLU, RoPE, “untied” embeddings), non-trainable RMSNorm parameters.
Width	[2048, 3072, 4096] (1024 for proxy model)
Depth	[16, 24, 32] (8 for proxy model)
Number of heads	[16, 24, 32] (8 for proxy model)
Head dimension	128
Total parameters	[1.07B, 3.12B, 6.98B]
Batch size	1024
Training steps	72000 (~ 300B tokens; 20000 for proxy model)
LR schedule	Cosine to 10%, 500 steps warm-up
Optimizer	AdamW ($\beta_1, \beta_2, \epsilon$) = (0.9, 0.95, 10^{-8})
Weight decay	2^{-13} , independent (Loshchilov & Hutter, 2019)
Dropout	0.0

Table 6: Large-scale training settings.

We use mixed-precision during training with optimizer states in FP32 that are sharded via ZeRO stage 1 (Rajbhandari et al., 2020). We retain the model weights in BF16 and apply our FP8 scheme as described in Section 3.2 to the tensors participating in matmul operations throughout the transformer block. All other tensors remain either in BF16 (embedding, readout layer, norm, activation function) or FP32 (Flash Attention (Dao et al., 2022)).

Each model was trained on several Nvidia A100 (80GB) or H100 GPUs, with all FP8 experiments conducted on the H100 chips utilizing their native FP8 support. For the FP8 operations we use PyTorch’s `torch._scaled_mm` function as a backbone.

B UNIT-SCALED OP DEFINITIONS

Table 7: Implementations of unit-scaled ops, building on Table A.2. from the Unit Scaling paper (Blake et al., 2023). These are considered part of u- μ P and should be used in the place of standard operations.

Op	Unit Scaling factors
$\text{matmul}(x, w) = xw$	$\alpha = \frac{1}{\sqrt{\text{fan-in}}}, \beta_x = \frac{1}{\sqrt{\text{fan-out}}}, \beta_w = \frac{1}{\sqrt{\text{batch-size}}}$
$\text{attention}(q, k, v) =$ $\text{softmax}(\alpha_{\text{attn}} d_{\text{head}}^{-1} (qk^\top) \odot c_{\text{mask}}) v$	$\alpha = \beta_q = \beta_k = \beta_v =$ $1 / \log_interpolate\left(\frac{1}{1 + \frac{4d_{\text{head}}}{\alpha_{\text{attn}}^2}}, 1, \sqrt{\frac{\log(s)}{s}}\right)$
$\text{gated_silu}(x_{\text{in}}, x_{\text{gate}}) =$ $x_{\text{in}} \odot x_{\text{gate}} \odot \text{sigmoid}(\alpha_{\text{ffn-act}} x_{\text{gate}})$	$\alpha = \beta_{x_{\text{in}}} = \beta_{x_{\text{gate}}} =$ $1 / \log_interpolate\left(\frac{1}{1 + \frac{1}{\alpha_{\text{ffn-act}}^2}}, \frac{1}{\sqrt{2}}, \frac{1}{2}\right)$
$\text{residual_add}(x_{\text{resid.}}, x_{\text{skip}}) =$ $a x_{\text{resid.}} + b x_{\text{skip}}$	$a = \frac{\tau}{\sqrt{\tau^2 + 1}}, b = \frac{1}{\sqrt{\tau^2 + 1}}$ (see G.2.2 for full details, inc. values for τ , which depends on α_{res} and $\alpha_{\text{res-attn-ratio}}$.)
$\text{softmax_xent}(x, t) =$ $\log_softmax(\alpha_{\text{loss-softmax}} x)_t$	$\alpha = 1, \beta = s / \sqrt{s - 1}$
RoPE(x)	$\alpha = \beta = 1$ (i.e. no scaling)
RMSNorm(x) (non-trainable, see (Lingle, 2024))	$\alpha = \beta = 1$ (i.e. no scaling)

The Unit Scaling paper provides scaling factors for various ops, in order to make them unit-scaled. However, these ops do not cover every case required for the Llama architecture used in our experiments, nor do they cover our updated residual layer implementation. To address this, in this section we outline a series of new unit-scaled ops for each of our required architectural features, as well as existing unit-scaled ops, as given in Table 7.

The presentation here is derived from that of the Unit Scaling Compendium given in (Blake et al., 2023, Table A.2). This makes reference to the factors $\alpha, \beta_1, \dots, \beta_k$. α is the output scaling factor in the forward pass, and β_i are the scaling factors for the gradient of the op’s inputs in the backward pass. For each op, a value or rule is provided for determining the required mult to ensure unit-scale. The correct value for these multipliers is derived by analyzing the scaling behavior of each op, given some reasonable distributional assumptions about the input and incoming gradient tensors (see Appendix E.2 for an example). Below we provide an in-depth overview of each new or modified unit-scaled op we introduce here.

Unit-scaled dot-product attention The Unit Scaling paper considers the attention layer scaling in terms of its separate components: the various matmul operations and the internal softmax. Linear operations are scaled using the standard rule, and the softmax scaling is given a $\alpha = \beta = s$ factor.

From an implementation perspective, the self-attention layer is more typically broken down into weight-matmuls and a fused scaled-dot-product attention operation. This is the case we handle here, accounting for three complicating factors not considered in the Unit Scaling paper:

1. As we use a decoder-style transformer in our experiments, our softmax operation has a causal mask applied to its input.

2. We follow the μP guidance of using $1/d_{\text{head}}$ scaling in our self-attention layer, rather than the usual $1/\sqrt{d_{\text{head}}}$.
3. We place a α_{attn} multiplier immediately before the softmax, which is an HP that users may tune.

As a result our dot-product attention takes the form:

$$\text{attention}(q, k, v) = \text{softmax}(\alpha_{\text{attn-softmax}} \cdot d_{\text{head}}^{-1} \cdot (q \cdot k^{\top}) \odot c_{\text{mask}}) \cdot v$$

The addition of an HP before the softmax introduces an additional challenge for Unit Scaling, as our scaling multipliers will need to account for this value when preserving unit scale.

This operation is sufficiently complex that we found an empirical model of its scale to be more accurate than any mathematically-derived rule (future work may consider justifying our model mathematically). We find that the scale of dot-product attention is approximately

$$\sigma(\text{attention}(q, k, v)) = \text{log_interpolate}\left(\frac{1}{1 + \frac{4d_{\text{head}}}{\alpha_{\text{attn}}^2}}, 1, \sqrt{\frac{\log(s)}{s}}\right)$$

where

$$\text{log_interpolate}(\alpha, b_{\text{upper}}, b_{\text{lower}}) = e^{\alpha \log(b_{\text{upper}}) + (1-\alpha) \log(b_{\text{lower}})}.$$

The corresponding scaling rule is therefore to divide by this factor in both the forward and backward pass, as outlined in Table 7.

SwiGLU FFN Llama uses a SwiGLU (Shazeer, 2020) layer for its FFN, which introduces two new operations for us to unit-scale: a SiLU (Yu & Su, 2019) (a.k.a. swish (Ramachandran et al., 2018)) operation and an element-wise multiplication. We take a similar approach to our dot-product attention, and consider unit-scaling the following fused operation:

$$\text{gated_silu}(x_{\text{in}}, x_{\text{gate}}) = x_{\text{in}} \odot x_{\text{gate}} \odot \text{sigmoid}(\alpha_{\text{ffn-act}} x_{\text{gate}})$$

For the surrounding weight-matmuls we follow the standard Unit Scaling rules.

Again, we use an empirical model of the scale of this op, which is surprisingly similar to the dot-product attention model:

$$\sigma(\text{gated_silu}(x_{\text{in}}, x_{\text{gate}})) = \text{log_interpolate}\left(\frac{1}{1 + \frac{1}{\alpha_{\text{ffn-act}}^2}}, \frac{1}{\sqrt{2}}, \frac{1}{2}\right),$$

dividing through by this factor to get our scaling rule.

Residual layers Our implementation of residual layers for u- μP is more complex than other operations, as adjustments are required to:

1. Make pre-norm residual networks support Unit Scaling (see Appendix F).
2. Introduce our new, principled residual HPs (see Appendix G).

Our residual layer scheme is presented in full in G.2.2. For readers interested in our justification for this, see the sections noted above.

We also follow the example of Unit Scaling and delay the application of our residual multiplier in the backward pass to the base of the branch (see (Blake et al., 2023), Figure 3c). This does not change the model, and enables unit-scale to be maintained on the residual branch regardless of the value of the multiplier.

RoPE embeddings We also require a unit-scaled implementation of Rotary Position Embeddings (RoPE (Su et al., 2024)), which are applied just before the scaled dot-product attention operation. As RoPE essentially consists of pair-wise rotations of elements by different degrees, we observe no meaningful scale-change as a result of its application, and hence leave it unchanged.

RMSNorm Following (Lingle, 2024) we opt to use a non-trainable version of RMSNorm (Zhang & Sennrich, 2019), in order to facilitate better transfer. As a result, we also leave this operation unchanged. Were a trainable RMSNorm to be used, the recipe would follow closely that of the LayerNorm presented in the Unit Scaling paper’s compendium.

Scale constraints One final, minor deviation from the scheme outlined in the Unit Scaling paper is the way in which we apply scale constraints (see their Section 5.2). The essence of scale constraints is that for perfect unit scaling, sometimes the ideal scale for the forward pass differs from those in the backward pass. In some special cases (e.g. at the ends of the network) the use of different scales can be valid, but in the general case a single scale must be agreed upon. The solution in the Unit Scaling paper is to use the geometric mean of the forward and backward scales.

We propose instead to simply use the forward scale over the backward scale(s) in these cases. We do so for the following reasons:

1. For these architectures we find empirically that where there is a disparity in ideal forward and backward scales, it is not large.
2. By taking the forward scale, we can ensure strict unit-scale in the forward pass.

The value of the latter point is in terms of what it means for the interpretation of our μP multiplier HPs. Consider the $\alpha_{\text{ffn-act}}$ multiplier; with strict unit scale we can say that the standard deviation of activations immediately before this multiplier is 1. Therefore the standard deviation immediately after is $\alpha_{\text{ffn-act}}$. As this multiplier is (by design) the last operation before the ffn activation function, we can say that the interpretation of $\alpha_{\text{ffn-act}}$ is simply to set the input standard deviation to the FFN’s activation function. Similar arguments can be made for other μP multiplier HPs. This interpretation only holds because we use the forward-scale in our constraints.

C THE CHALLENGES WITH μP IN PRACTICE

C.1 NOT ALL TRAINING SETUPS GIVE μT RANSFER

Lingle (2024) shows that directly applying μP to a decoder LM fails to provide LR transfer across width. Given that the primary use of μP in the literature has been LM training of this kind, this result suggests a significant limitation. How do we reconcile this with the strong LR transfer across width shown for language models in Tensor Programs V?

We answer this in Figure 19. The first training setup (a) is aligned with that used in Tensor Programs V (their Figure 4). There are several atypical aspects to their training setup, primarily the use of a constant LR schedule and a high number of epochs; we outline the precise differences between setup (a) and (b) in Table 8. This overfitting regime makes validation loss unusable, and transfer misleadingly good. When we remove these and shift to a standard Llama training setup (b), optimal HPs begin to drift with width (see Figure 21 for an ablation of individual changes). This confirms Lingle’s findings that standard μP is in fact a poor fit for modern LM training. We fix this (c) by the removal of parameters from LayerNorms/RMSNorms, as suggested by Lingle, and the introduction of *independent* weight decay for AdamW, as suggested by Wortsman et al. (2023)⁴ (see Wang & Aitchison (2024) for further analysis). With these changes adopted, we recover the strong transfer shown in Tensor Programs V’s experiments. Each change is evaluated independently in Figure 20, which shows that the dominant effect is a narrowing of the learning basin due to non-parametric RMSNorms, leading to better learning rate transfer.

C.2 IT’S NOT CLEAR WHICH HYPERPARAMETERS TO SWEEP

The problem of selecting HPs to sweep can be framed as choosing a subset of the per-tensor $\alpha_W, \sigma_W, \eta_W$ HPs outlined in Section 2.1, and grouping across/within layers. As shown in Table 9, μT ransfer experiments in the literature have done this in a variety of ways. Practitioners have not

³ As in other work, we use μP as a shorthand for the method outlined in Tensor Programs V, including μT ransfer. Strictly speaking, μP ought only to refer to the parametrization outlined in Tensor Programs IV. ⁴ Lingle suggests independent weight decay is unstable, but we find it to be more so than Adam or standard AdamW.

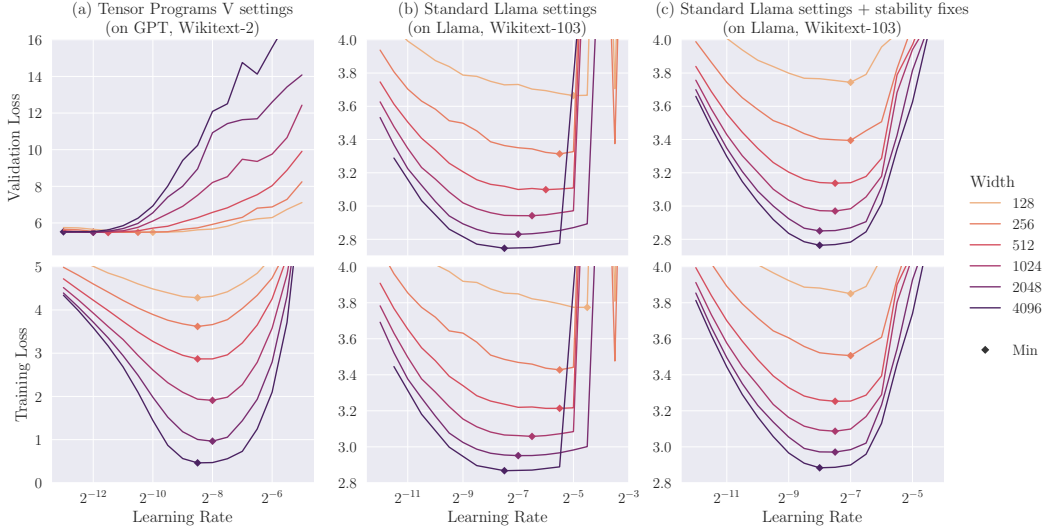


Figure 19: Effective μ Transfer does not hold across all training setups. **(a)** We show strong transfer for the unrealistic setup used in Tensor Programs V (too many epochs; constant LR). **(b)** Moving to a more standard Llama training setup, transfer breaks down. **(c)** This is restored by the introduction of two improvements to transfer stability: non-parametric norms and independent weight decay.

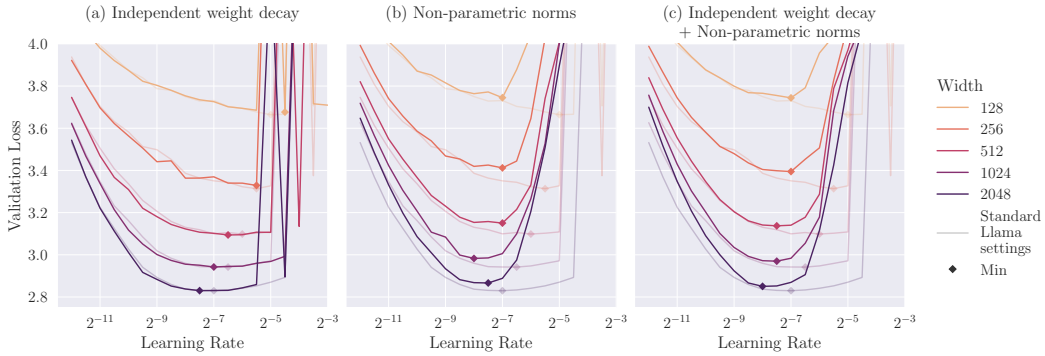


Figure 20: The effect of the individual transfer stability fixes from Figure 19. **(a)** In this setting switching from non-independent to independent weight decay has only a minor effect, though Wortsman et al. (2023) Figure 6 suggests it may be highly valuable in other settings. **(b)** Non-parametric norms give a narrower learning rate basin, leading to better transfer. **(c)** The combination of these, for comparison, matching Figure 19 (c).

Feature	Tensor Programs V	Standard Llama
Dataset	wikitext-2	wikitext-103
Vocab Size	33278	32000
Nsteps	10000	8192
Batch Size	20	64
Optimizer	adam	adamw
LR Schedule	constant	cosine
Weight Decay	0	0.00012
Positional Encoding	absolute	rotary
Norm	layer_norm	rms_norm
Dropout	0.2	0
NLayers	2	4
Use Gated FFN	False	True
Activation FN	relu	swish
FFN Ratio	4	2.75
Final Norm	False	True
Base Depth	1	4
Zero Init Readout	True	False

Table 8: Comparison of Tensor Programs V’s standard settings (as best we can tell) and our Standard Llama setup, corresponding to (a) and (b) in Figure 19.

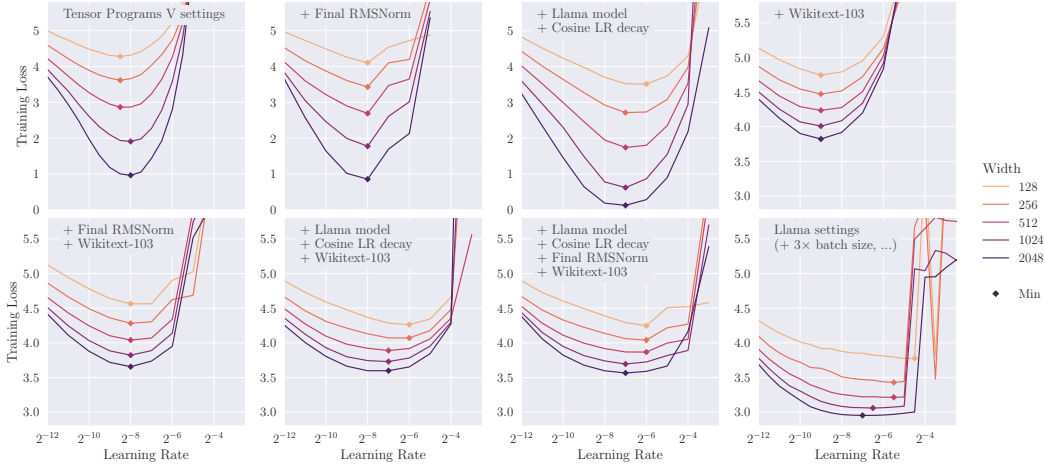


Figure 21: An ablation of the more standard Llama training settings against the Tensor Programs V settings from Figure 19. This shows that the flat basins with poor transfer are not due to a single change, but the combination of a larger dataset (training < 1 epoch) and the stronger Llama model are largely responsible. Note that ‘Llama model’ here indicates a group of changes: rms norm, rotary embeddings & swiglu FFN.

justified these choices, appearing to rely on a mixture of precedent and intuition. We outline two major downsides to the lack of a principled approach.

Firstly, not all groupings of HPs are suitable. Consider the commonly-used global σ_{init} HP. At initialization the activations going into the FFN swish function have $\text{std}(x_{\text{swish}}) \propto \sigma_{W_{\text{gate}}}$, whereas the self-attention softmax activations have $\text{std}(x_{\text{attn}}) \propto \sigma_{W_Q} \sigma_{W_K}$. A global σ HP thus has a linear effect on the FFN and a quadratic effect on attention, suggesting that this grouping may be unideal.

Secondly, not all HPs are independent of one another. The key example of this is the interaction between σ_W and η_W . The relative size of a weight update is determined by the ratio η_W / σ_W , not by either HP individually. Because of this, the optimal values for σ and η depend on each other, which we demonstrate empirically in Section 4.1. This can make the problem of HP search much harder, and may be why hundreds of random-search runs have been required for sweeps in the literature.

C.3 BASE SHAPE COMPLICATES USAGE

Most practitioners are unlikely to require alignment with an SP model, in which case it is unclear what base-width (and base-depth) should be used. The literature has aligned on a standard base-width of 256 (see Table 9), but this appears to lack a principled motivation—though the fact that they are not dropped entirely suggests they may be beneficial under u- μ P.

Implementing base-shape HPs (see Equation (3)) can also add complications from an engineering perspective. The proposed implementation in the `mup` library (Microsoft, 2024) reflects this, requiring an extra ‘base’ model to be created and the original model to be re-initialized. This can interact awkwardly with other model-transforms for features like quantization, compilation, etc:

```
import mup

proxy_model = MupModel(d_model=128, ...) # proxy width
base_model = MupModel(d_model=256, ...) # base width
mup.set_base_shapes(proxy_model, base_model) # re-init proxy_model
```

C.4 μ P APPEARS TO STRUGGLE WITH LOW-PRECISION

Finally, we note an interesting contradiction observed in the relationship between μ P and low-precision. One of the stated aims for μ P is that its activations have $\Theta(1)$ -sized coordinates in the limit (Yang et al., 2022, Desiderata J.1). This desideratum is specifically given in order that values can be represented using finite-range floating-point numbers (Yang & Hu, 2021, Section 3). Yet despite numerical stability being central to the theory underlying μ P, this is not leveraged to ensure that μ P models can *actually* be trained in low-precision. Indeed, for the LLM runs in Tensor Programs V the SP model trains successfully in FP16, while the μ P model diverges (attributed to underflow of gradients). We remedy this with u- μ P.

D A GUIDE TO USING U- μ P

We bring together our u- μ P scheme presented in Section 3 to form a simple recipe for applying it to a model. The u- μ P scheme is designed and validated on a Llama-style architecture, so it may not be applicable or effective on other models, particularly those with substantially different architectures. Exploring this question is an important avenue for future work.

Before applying our scheme, users are encouraged to apply the following pre-requisites to their training setup, based on our analysis of effective μ Transfer in Appendix C.1:

- Remove trainable parameters from normalization layers
- Use the *independent* form of AdamW
- Ensure training is in the under-fitting regime (i.e. avoid excessive data repetition)

Having done this, our recipe for using u- μ P is as follows:

1. **Replace operations & optimizers with u- μ P versions:** Each operation should be replaced by a unit-scaled version (these wrap the existing operations, with added static scales in the forward and backward passes). We have pre-calculated scales for common operations in Appendix B. Parameters should be initialized with unit variance, and Adam(W) adjusted to use the scaling rules defined in Section 3.4 (we refer to the optimizer as Adam in this section, but AdamW should be used if weight decay is required. Other optimizer scaling rules can be determined by the same process we outline).
2. **Choose a set of HPs to sweep:** From the set of HPs outlined in Table 3, select those to be swept. We recommend the extended set, though a basic LR sweep can be effective.
3. **Decide on proxy model config:** The cost of proxy model training should be such that the sweeping process is much less than target model training, while still being as representative as possible. We base our recommendations on the results in Figure 4. In general, width is

the most reliable feature to transfer. Training steps and batch size also give good transfer, so moderate changes here are permissible. Depth is the least reliable feature for transfer, so we only recommend modest changes in depth. We keep the number of warmup steps constant, but always decay to the same final LR when varying the number of steps.

4. **Perform independent HP search:** Following the process outlined in Section 4.1 and Appendix A.5.
5. **Train the target model:** This can be done in FP8 simply by placing casts on matmul inputs (though for our large-scale experiments we found the scales of two operations drifted enough over time that some lightweight dynamic re-scaling was required).

E ADDITIONAL BACKGROUND MATERIAL

E.1 THE MAXIMAL UPDATE PARAMETRIZATION

Theoretical background We do not cover the theory underpinning μP in this paper, presenting only its resulting scaling rules (Table 1). For readers interested in this theory, the extensive Tensor Programs series (Yang, 2019; 2020a; Yang & Littwin, 2021; Yang, 2020b; Yang & Littwin, 2023) builds up a framework from which μP is derived (Yang & Hu, 2021). For those requiring a more accessible introduction, Yang et al. (2023a) show that μP can be derived in a simpler and more general way by placing a spectral scaling condition on the norm of weights and their updates.

Approaches to HP sweeping in the literature Table 9 outlines the ways in which users of μP in the literature have approached HP sweeping. These all follow the approach used in Tensor Programs V of a random sweep, sampling combinations from the joint space of all HPs. The authors of Tensor Programs V note that other more complex methods may be more efficient, but these are considered beyond the scope of their work and have not been used widely. A Bayesian search method was used for the development of MiniCPM (Hu et al., 2024), but the authors give no further details—as they use 400 runs in their sweep it is not clear that this approach makes HP search easier.

Table 9: Sweeping configurations used for a selection of μP models from the literature. The sweeping process is similar across models, the only differences being the choice of discrete or continuous distributions and their ranges. Model references: T.P.V WMT14 (Yang et al., 2022), T.P.V BERT_{large} (Yang et al., 2022), T.P.V GPT-3 (Yang et al., 2022), MiniCPM (Hu et al., 2024), Cerebras-GPT (Dey et al., 2023a), S μ Par (Dey et al., 2024b).

Model	proxy/target tokens used	proxy/target model size	sweep size	base width	HPs swept
T.P.V WMT14	100%	7.1%	64		$\eta, \alpha_{\text{out}}, \alpha_{\text{attn}}$
T.P.V BERT _{large}	10%	3.7%	256	?	$\eta, \eta_{\text{emb}}, \alpha_{\text{out}}, \alpha_{\text{attn}}, \alpha_{\text{LN}}, \alpha_{\text{bias}}$
T.P.V GPT-3	1.3%	0.6%	350		$\eta, \sigma, \alpha_{\text{emb}}, \alpha_{\text{out}}, \alpha_{\text{attn}}, \alpha_{\text{pos}}$
MiniCPM	0.008%	0.45%	400	256	$\eta, \sigma, \alpha_{\text{emb}}, \alpha_{\text{residual}}$
Cerebras-GPT	1.1%	1.5%	200	256	$\eta, \sigma, \alpha_{\text{emb}}$
S μ Par	6.6%	6.4%	350	256	$\eta, \sigma, \alpha_{\text{emb}}$

E.2 UNIT SCALING

An example: the unit-scaled matmul op Here we outline the procedure for calculating the scaling factor of a matmul op, which practitioners can use as a guide for scaling new ops that we do not cover in this paper (see Appendix B).

There are two potential approaches here. The first is to derive scaling factors from an analysis of an op’s dynamics. Specifically, given the assumption of unit-scaled inputs, the appropriate scaling factor is the reciprocal of the expected output scale. For a basic matrix-matrix matmul we have,

$$\text{matmul}(X, W) = XW, \quad X \in \mathbb{R}^{d_{\text{batch}} \times d_{\text{fan-in}}}, \quad W \in \mathbb{R}^{d_{\text{fan-in}} \times d_{\text{fan-out}}},$$

where weights and activations are sampled i.i.d. from a centered Gaussian:

$$X_{ij} \sim \mathcal{N}(0, \sigma_X^2), \quad W_{jk} \sim \mathcal{N}(0, \sigma_W^2).$$

From this we can derive the expected output scale (i.e. $\sigma(\text{matmul})$):

$$\text{matmul}(X, W)_{ik} = \sum_{j=1}^{d_{\text{fan-in}}} X_{ij} W_{jk},$$

$$\sigma(\text{matmul}(X, W)_{ik}) = \sqrt{d_{\text{fan-in}}} \sigma_W \sigma_X.$$

Under Unit Scaling we have $\sigma_W = \sigma_X = 1$, and hence the scaling factor required to ensure a unit-scaled output is $1/\sqrt{d_{\text{fan-in}}}$. This gives our final unit-scaled matmul:

$$\text{u-matmul}(X, W) = \text{matmul}(X, W) / \sqrt{d_{\text{fan-in}}}$$

The distributional assumptions made here hold at initialization, but do not over training. A more precise model for the asymptotic behavior of neural networks under training is given by the Tensor Programs framework, but for the purposes of numerics this precise treatment of scale at initialization appears to be sufficient.

The second, less ideal approach to calculating scaling factors is to use experimentation to infer this relationship empirically. In this case, one would sample random initializations and compute the output scale over a range of $d_{\text{fan-in}}$ values (or whatever HPs one expects the output scale to depend on), fitting a curve to the observed data.

Applying unit scaling To apply Unit Scaling to a model and train in low-precision, the following steps are required:

1. Scale parameter initializations to have zero-mean and unit variance.
2. Replace operations with their unit-scaled equivalents (including and especially the loss, matmuls and residual-adds).
3. *Constrain* the scales of operations which are required to have the same forward and backward factors.
4. Place a simple `.to(fp8)` cast on the inputs to matmuls.

Step 3 relates to the problem of conflicting scales in the forward and backward passes. A single linear layer in a differentiated model requires 3 matmul ops in the forward and backward passes, each requiring a different scaling factor ($\frac{1}{\sqrt{d_{\text{fan-in}}}}$, $\frac{1}{\sqrt{d_{\text{fan-out}}}}$, $\frac{1}{\sqrt{d_{\text{batch-size}}}}$). However, using these directly would give invalid gradients. The compromise here is that the activations and activation gradients have their scaling factors *constrained* such that they are equal (the Unit Scaling paper recommends taking the geometric mean; we modify this for u- μ P in Appendix B to simply use the forward scale everywhere). Weight gradients can still be given their own scaling factor due to the *cut-edge rule* (as explained in Appendix H).

Step 4 reflects the key benefit of Unit Scaling. Unlike other methods it changes the learning dynamics of a model, but the advantage is that unit-scaled models then ‘naturally’ generate well-scaled tensors. This means that low-precision arithmetic ideally becomes as simple as placing a cast operation before matmuls as outlined.

F UNIT-SCALED PRE-NORM RESIDUAL LAYERS

The popular pre-norm residual network architecture is simple to implement, but problematic to combine with Unit Scaling. It exhibits scale-growth in the skip-stream at initialization, due to the repeated addition of residual connections without subsequent normalization. Here we present a surprising and useful finding: that for any pre-norm model there exists a mathematically-equivalent model where this scale-growth is eliminated, through the careful re-scaling of residual connections.

Note that this section focuses on applying Unit Scaling to *standard* pre-norm models. Only once we have addressed this problem are we able to do the same for u- μ P models, as shown in Appendix G.2. Readers only interested in our final u- μ P residual implementation may skip ahead to Appendix G.2.2.

F.1 SCALE GROWTH IN PRE-NORM RESIDUAL NETWORKS

Let's consider a pre-norm residual network of depth L :

$$R_0(x) = r_0 x, \quad (4)$$

$$R_l(x) = r_l f_l(R_{l-1}(x)) + R_{l-1}(x), \quad l = 1, \dots, L \quad (5)$$

$$R_{L+1}(x) = f_{L+1}(R_L(x)) \quad (6)$$

with embedding multiplier r_0 and residual branch multipliers r_l for $l = 1, \dots, L$. To satisfy pre-norm, all f_l are zero-homogeneous functions, i.e. $f_l(\lambda x) = f_l(x)$.

The scale of the skip-stream at initialization as a result of Equation (5) is

$$\sigma(R_l) = \sqrt{r_l^2 \sigma(f_l)^2 + \sigma(R_{l-1})^2} > \sigma(R_{l-1}), \quad l = 1, \dots, L \quad (7)$$

assuming $r_l^2 \sigma(f_l)^2 > 0$. This shows that scale inevitably grows with the addition of each residual layer.

This scale-growth is clearly incompatible with unit scaling, which aims for $\sigma(R_l) = 1$ for all $l = 0, \dots, L + 1$. In the following we present an elegant solution to this problem making use of a symmetry transformation available in pre-norm residual architectures.

F.2 RESIDUAL SYMMETRY IN PRE-NORM ARCHITECTURES

To resolve the problem of scale shift in residual networks demonstrated by Equation (7), we try a slightly more general ansatz:

$$\hat{R}_0(x) = x, \quad (8)$$

$$\hat{R}_l(x) = a_l f_l(\hat{R}_{l-1}(x)) + b_l \hat{R}_{l-1}(x), \quad (9)$$

$$\hat{R}_{L+1}(x) = f_{L+1}(\hat{R}_L(x)) \quad (10)$$

with coefficients a_l, b_l . We want to choose these coefficients so that the outputs of \hat{R}_l are unit-scaled if the outputs f_l, \hat{R}_{l-1} are. A similar calculation as in Equation (7) leads to the sufficient condition

$$a_l^2 + b_l^2 = 1, \quad (11)$$

which can be easily satisfied. Having restored Unit Scale, we are faced with another issue. It seems that Equations (8) to (10) describe a different network than Equations (4) to (6), whereas ideally the relation from input to final output should be unchanged when converting the network to Unit Scaling.

Note that the coefficients a_l, b_l are not uniquely defined yet, so our mathematical intuition tells us that we should find an additional constraint to get a unique solution. To find this constraint, let us consider our original residual network in Equations (4) to (6) and analyze how the variance propagates through the network if we assume all the f_l satisfy Unit Scaling and $\sigma(x) = 1$. Let σ_{l-1}^2 denote the variance of R_{l-1} . Then a simple inductive calculation shows that

$$\sigma_{l-1}^2 = \sum_{i=0}^{l-1} r_i^2.$$

By Equation (5) the output of R_l adds a quantity of scale r_l from the residual connection and a quantity of scale σ_{l-1} from the skip connection. Intuitively, the *ratio* of these scales should be more important for the overall network dynamics than their absolute values. Thus our constraint becomes preserving the ratio of scales from the original model, through our choice of a_l, b_l :

$$\frac{a_l}{b_l} = \frac{\sigma(r_l f_l)}{\sigma_{l-1}} = \frac{r_l}{\sqrt{\sum_{i=0}^{l-1} r_i^2}} =: \tau_l,$$

which, recalling Equation (11), (up to sign) uniquely defines our multipliers a_l, b_l as

$$a_l = \frac{\tau_l}{\sqrt{\tau_l^2 + 1}}, \quad b_l = \frac{1}{\sqrt{\tau_l^2 + 1}} \quad (12)$$

In summary, we propose the modified residual network

$$\hat{R}_0(x) = x, \quad (13)$$

$$\hat{R}_l(x) = \frac{\tau_l}{\sqrt{\tau_l^2 + 1}} f_l(\hat{R}_{l-1}(x)) + \frac{1}{\sqrt{\tau_l^2 + 1}} \hat{R}_{l-1}(x), \quad (14)$$

$$\hat{R}_{L+1}(x) = f_{L+1}(\hat{R}_L(x)), \quad (15)$$

$$\tau_l^2 = \frac{r_l^2}{\sum_{i=0}^{l-1} r_i^2}. \quad (16)$$

Our main result of this section is that this network is indeed mathematically equivalent to the network defined in Equations (4) to (6), under a simple additional structural assumption:

Lemma F.1. *Consider R_l , \hat{R}_l defined as in Equations (5) and (14) respectively. Then $\hat{R}_l = R_l / \sqrt{\sum_{i=0}^l r_i^2}$ for all $l = 0, \dots, L$.*

Remarkably, this result does not assume the individual network operations f_l actually satisfy Unit Scaling. It is purely a consequence of the pre-norm residual structure. However, only under Unit Scaling can the factors τ_l be interpreted as the ratio of scales between skip and residual branch.

As a consequence of the lemma, the final residual output $R_L(x)$ is the same as in our original network up to a fixed multiplier. Due to the zero-homogeneity of the final output function f_{L+1} this gives $\hat{R}_{L+1} = f_{L+1} \left(R_L(x) / \sqrt{\sum_{i=0}^L r_i^2} \right) = f_{L+1}(R_L(x)) = R_{L+1}$, proving the mathematical equivalence of our residual scheme. Modern LLM architectures like Llama (Touvron et al., 2023a) are pre-norm residual networks of this kind. Hence they admit a faithful unit-scaled reparametrization.

F.3 PROOF OF LEMMA F.1

Proof. This is proved by induction. For the base-case $l = 1$, we have $\tau_1 = r_1/r_0$, giving

$$\begin{aligned} \hat{R}_1(x) &= \frac{\tau_1}{\sqrt{\tau_1^2 + 1}} f_1(x) + \frac{1}{\sqrt{\tau_1^2 + 1}} x \\ &= (r_1 f_1(x) + r_0 x) / \sqrt{r_0^2 + r_1^2} \\ &= R_1 / \sqrt{r_0^2 + r_1^2}. \end{aligned}$$

Then if the statement holds for $l - 1$ we have

$$\begin{aligned} \hat{R}_l(x) &= \frac{\tau_l}{\sqrt{\tau_l^2 + 1}} f_l(\hat{R}_{l-1}(x)) + \frac{1}{\sqrt{\tau_l^2 + 1}} \hat{R}_{l-1}(x) \\ &= \frac{r_l}{\sqrt{\sum_{i=0}^l r_i^2}} f_l(\hat{R}_{l-1}(x)) + \frac{\sqrt{\sum_{i=0}^{l-1} r_i^2}}{\sqrt{\sum_{i=0}^l r_i^2}} \hat{R}_{l-1}(x) \\ &= \left(r_l f_l(\hat{R}_{l-1}(x)) + \sqrt{\sum_{i=0}^{l-1} r_i^2} \hat{R}_{l-1}(x) \right) / \sqrt{\sum_{i=0}^l r_i^2} \\ &= \left(r_l f_l(R_{l-1}(x)) + \sqrt{\sum_{i=0}^{l-1} r_i^2} \frac{R_{l-1}(x)}{\sqrt{\sum_{i=0}^{l-1} r_i^2}} \right) / \sqrt{\sum_{i=0}^l r_i^2} \\ &= (r_l f_l(R_{l-1}(x)) + R_{l-1}(x)) / \sqrt{\sum_{i=0}^l r_i^2} \\ &= R_l(x) / \sqrt{\sum_{i=0}^l r_i^2} \end{aligned}$$

□

F.4 UNIT SCALING FOR TRANSFORMER RESIDUALS

The above scheme describes Unit Scaling for arbitrary pre-norm residual networks. We now apply it to the case of pre-norm transformer residual layers.

We can describe a transformer in terms of the residual network given in Equations (4) to (6). Our f_l functions alternate between self-attention layers and feed-forward layers. Implementations differ in the handling of how residual multipliers r_l correspond to HPs. In many cases practitioners simply ignore these r_l , but for the sake of expressivity we assume the two types of residual layer each have their own HP, as well as the embedding. In other words,

$$r_l = \begin{cases} \alpha_{\text{emb}} & l = 0 \\ \alpha_{\text{attn-residual}} & l \text{ is odd} \\ \alpha_{\text{ffn-residual}} & l \text{ is even, and } l > 0. \end{cases}$$

To convert this to a Unit Scaled network we apply Equations (13) to (16), from which can derive the following closed-form expression for τ_l :

$$\tau_l^2 = \begin{cases} \frac{\alpha_{\text{attn-residual}}^2}{\alpha_{\text{emb}}^2 + \ell \alpha_{\text{attn-residual}}^2 + \ell \alpha_{\text{ffn-residual}}^2} & l \text{ is odd} \\ \frac{\alpha_{\text{ffn-residual}}^2}{\alpha_{\text{emb}}^2 + (\ell + 1) \alpha_{\text{attn-residual}}^2 + \ell \alpha_{\text{ffn-residual}}^2} & l \text{ is even.} \end{cases}$$

where $\ell = \lfloor \frac{l-1}{2} \rfloor$.

This gives us a unit-scaled pre-norm residual implementation for a *standard* transformer, which is mathematically equivalent to a non-unit-scaled version. In the next section we augment this by adding in two HPs, in a carefully-designed manner that satisfies our criteria for u- μ P HPs, giving us our full residual implementation.

G JUSTIFYING THE U- μ P HYPERPARAMETER SCHEME

Here we justify our particular choice of u- μ P HP, as given in Table 3 (with their placement defined in Table 7). We discuss this topic briefly in Section 3.3, stating that all our HPs (excepting the LR) are α HPs, and under u- μ P they are now associated with operations instead of weights. All operations have an α HPs, unless they are unary and k -homogeneous for $k \geq 0$.

We begin this section by explaining why we apply this rule to the model and how it results in three of our u- μ P HPs. We then consider how best to hyperparametrize our residual layers, building on our criteria for HPs given in Section 3.3 and the unit-scaled pre-norm residual scheme in Appendix F.

G.1 MULTIPLIERS FOR NON-HOMOGENEOUS OPS: $\alpha_{\text{attn-softmax}}$, $\alpha_{\text{ffn-act}}$, $\alpha_{\text{loss-softmax}}$

In this section we derive the rest of our u- μ P multipliers. We want to identify the minimal set that can still express all different choices of pre-op scales in the model. The crucial observation is that every pre-scale multiplier α of a unary operation $h \mapsto f(\alpha h)$ can be propagated through the network if f is k -homogeneous for some $k > 0$, i.e. $f(\alpha x) = \alpha^k f(x)$, leaving the model and its optimization unchanged. We can iterate this along the computational path until either the next operation is non-homogeneous, non-unary (we are at the end of a residual path), or the next operation is 0-homogeneous (e.g. a norm).

In the first case the accumulated scales are absorbed in the pre-op scale of the non-homogeneous operation (where we introduce a multiplier), in the second case they are absorbed in the residual addition for that branch (where we again introduce a multiplier), and in the final case the scale disappears (so we start over). We now go through the Llama forward computation and follow this paradigm to identify our multipliers in Table 10.

Table 10: A walkthrough of the Llama architecture, showing how our $\alpha_{\text{attn-softmax}}$, $\alpha_{\text{ffn-act}}$ and $\alpha_{\text{loss-softmax}}$ multipliers are derived via an analysis of scale-propagation.

Op	Scale propagation behavior
Embedding	We show in Appendix G.2.1 that the embedding multiplier can be absorbed in the residual multipliers, meaning one is not required here.
Attention RMSNorm	This operation is 0-homogeneous and thus we start over.
Query & key projection	Both are linear, meaning their scale is propagated. Multipliers are therefore not required.
Query-key matmul	Again linear. As query & key are both generated from the same input, this operation is 2-homogeneous wrt. that input. Hence it also propagates scale.
Softmax	The softmax operation is non-homogeneous. Thus the pre-op scale of the softmax becomes our first multiplier: $\alpha_{\text{attn-softmax}}$.
Value	The value layer is linear and hence propagates scale.
Softmax-value matmul	Again linear and hence propagates scale.
Attention projection	This operation is linear and lies at the end of the attention residual path. Hence there are no more multipliers in the attention block.
Residual add	This operation is non-unary and hence receives our second (and third) multipliers: α_{res} , $\alpha_{\text{res-attn-ratio}}$. The manner and motivation for using two multipliers here is justified in the next section.
FFN RMSNorm	This operation is 0-homogeneous and thus we start over.
FFN input scale	The input layer is linear, hence it propagates scale.
Sigmoid input	This function is non-homogeneous and thus we have our fourth multiplier: $\alpha_{\text{ffn-act}}$.
SiLU weight	This layer is also linear and propagates scale.
Product	The entry-wise multiplication of the outputs of sigmoid, input layer and SiLU weight is homogeneous and thus propagates scale.
FFN output	This layer is linear and at the end of the residual path. Hence there are no more multipliers in the FFN residual block.
Residual add	See above.
Output RMSNorm	This operation is 0-homogeneous and thus we start over.
Output head	This layer is linear, hence it propagates scale.
Loss	The cross-entropy loss is non-homogeneous and leads to our final multiplier: $\alpha_{\text{loss-softmax}}$.

G.2 RESIDUAL BRANCH MULTIPLIERS: α_{res} , $\alpha_{\text{res-attn-ratio}}$

In this section we derive our two u- μ P residual HPs. We start with the basic, non-unit scaled model we began with in the previous section, outlined in Equations (4) to (6). We described a set of α_{emb} , $\alpha_{\text{attn-residual}}$, $\alpha_{\text{ffn-residual}}$ HPs associated with this model in Appendix F.4. However these HPs poorly satisfy our cardinality, independence and interpretability criteria from Section 3.3, so in the Appendix G.2.1 we present a re-parametrization of these HPs designed to better satisfy these points. In Appendix G.2.2 we then combine these HPs with the final unit-scaled pre-norm residual scheme we derived in Appendix F, resulting in our complete u- μ P residual scheme.

G.2.1 IMPROVED HYPERPARAMETERS FOR TRANSFORMER RESIDUALS

To avoid cluttered notation, in this section we rename

$$\begin{aligned}\alpha_{\text{res}} &= \alpha_r, & \alpha_{\text{res-attn-ratio}} &= \alpha_\rho \\ \alpha_{\text{emb}} &= \alpha_e, & \alpha_{\text{attn-residual}} &= \alpha_a & \alpha_{\text{ffn-residual}} &= \alpha_f.\end{aligned}$$

To make the presentation more clear, we derive our new HPs using the standard residual scheme from Equations (4) to (6). For the actual unit scaled implementation one needs to transform the multipliers following Equations (13) to (16), which we do in Section G.2.2.

To facilitate our analysis, we can view the transformer residual output as the sum of three terms:

$$\begin{aligned}R_L &= R_L^{(e)} + R_L^{(a)} + R_L^{(f)}, \\ R_L^{(e)} &:= \alpha_e x, \\ R_L^{(a)} &:= \sum_{l=1}^{L/2} \frac{\alpha_a}{\sqrt{L/2}} f_{2l-1}(R_{2l-1}(x)), \\ R_L^{(f)} &:= \sum_{l=1}^{L/2} \frac{\alpha_f}{\sqrt{L/2}} f_{2l}(R_{2l}(x)),\end{aligned}$$

and define the average residual scale,

$$\sigma(R_L^{(a,f)})^2 := \frac{\sigma(R_L^{(a)})^2 + \sigma(R_L^{(f)})^2}{2}.$$

Note that we have added in the depth- μ P multipliers here, though a similar analysis can be performed for non-depth- μ P models. As above, f_l functions alternate between self-attention layers and feed-forward layers.

With respect to our interpretability criterion, we propose two new multipliers that correspond to dynamics in the network which we suggest are important to control at initialization. The first is the ratio of the average scale of the residuals' contributions to those of the embedding, $\alpha_r = \sigma(R_L^{(a,f)})/\sigma(R_L^{(e)})$. The second is the ratio of the scale of the attention-residuals' contributions to those of the feed-forward-residuals, $\alpha_\rho = \sigma(R_L^{(a)})/\sigma(R_L^{(f)})$. Not only do these two ratios control key dynamics of our model, but we can use them to replace our existing $(\alpha_e, \alpha_a, \alpha_f)$ multipliers.

Let us first examine these two quantities for a standard (non-unit-scaled model). Residual functions of the same kind have the same expected output scale at initialization in pre-norm networks, meaning we can denote the output scale $\sigma(f_l(R_l))$ of all self-attention functions as σ_a , and of all feed-forward functions as σ_f . We thus have the following scales at the output:

$$\begin{aligned}\sigma(R_L^{(e)}) &= \alpha_e \sigma(x), \\ \sigma(R_L^{(a)}) &= \frac{\alpha_a}{\sqrt{L/2}} \sigma\left(\sum_{i=1}^{L/2} f_{2i-1}(R_{2i-1})\right) = \alpha_a \sigma_a, \\ \sigma(R_L^{(f)}) &= \frac{\alpha_f}{\sqrt{L/2}} \sigma\left(\sum_{i=1}^{L/2} f_{2i}(R_{2i})\right) = \alpha_f \sigma_f, \\ \sigma(R_L^{(a,f)}) &= \sqrt{\frac{(\alpha_a \sigma_a)^2 + (\alpha_f \sigma_f)^2}{2}}.\end{aligned}$$

Recalling our definitions of α_r, α_ρ above, this gives us:

$$\begin{aligned}\alpha_\rho &= \frac{\alpha_a \sigma_a}{\alpha_f \sigma_f}, \\ \alpha_r &= \sqrt{\frac{(\alpha_a \sigma_a)^2 + (\alpha_f \sigma_f)^2}{2(\alpha_e \sigma(x))^2}}, \\ &= \sqrt{\frac{\alpha_\rho^2 + 1}{2}} \frac{\sigma_f}{\sigma(x)} \frac{\alpha_f}{\alpha_e}.\end{aligned}$$

The original α_a, α_f multipliers can then be written in terms of α_r, α_ρ :

$$\begin{aligned}\alpha_a &= \alpha_\rho \alpha_f \frac{\sigma_f}{\sigma_a} \\ \alpha_f &= \alpha_r \alpha_e \frac{\sigma(x)}{\sigma_f} \sqrt{\frac{2}{\alpha_\rho^2 + 1}}\end{aligned}$$

We have replaced two of the three original multipliers, but still have a dependence on α_e here in our expressions for α_f and $R_L^{(e)}$, which we now remove by dividing it out of our residual branches and embedding. We use the hat ($\hat{\cdot}$) symbol to denote terms that have been divided-through by α_e . This new system of equations is equivalent to our old one thanks to the zero-homogeneity of the final post-residual layer:

$$\begin{aligned}R_{L+1}(x) &= f_{L+1}(R_L^{(e)} + R_L^{(a)} + R_L^{(f)}) \\ &= f_{L+1}((R_L^{(e)} + R_L^{(a)} + R_L^{(f)})/\alpha_e) \\ &= f_{L+1}(\hat{R}_L^{(e)} + \hat{R}_L^{(a)} + \hat{R}_L^{(f)})\end{aligned}$$

This gives $\hat{R}_L^{(e)} = \alpha_e x / \alpha_e = x$, removing our first occurrence of α_e . Following the division through $\hat{R}_L^{(a)}$ and $\hat{R}_L^{(f)}$, we obtain:

$$\begin{aligned}\hat{R}_L^{(a)} &:= \sum_{l=1}^{L/2} \frac{\hat{\alpha}_a}{\sqrt{L/2}} f_{2l-1}(R_{2l-1}), \\ \hat{R}_L^{(f)} &:= \sum_{l=1}^{L/2} \frac{\hat{\alpha}_f}{\sqrt{L/2}} f_{2l}(R_{2l}), \\ \hat{\alpha}_a &= \alpha_\rho \hat{\alpha}_f \frac{\sigma_f}{\sigma_a}, \\ \hat{\alpha}_f &= \alpha_r \frac{\sigma(x)}{\sigma_f} \sqrt{\frac{2}{\alpha_\rho^2 + 1}}.\end{aligned}$$

This system of equations is the same as the original, but with the two α_e terms dropped, meaning our model's multipliers can be expressed in terms of only α_r and α_ρ . Using the above equations, any pair of values for (α_r, α_ρ) can be translated back into an equivalent set of values for $(\alpha_e, \alpha_a, \alpha_f)$ such that the output $R_{L+1}(x)$ is the same, meaning that our multipliers are no less expressive than the original set. This satisfies our desired criteria of minimizing the number of multipliers while maintaining expressivity.

We can simplify further in the case of unit-scaled models, which are designed such that $\sigma(x), \sigma_a, \sigma_f$ are all 1 at initialization. In this case our re-parametrization becomes:

$$\hat{\alpha}_a = \alpha_\rho \hat{\alpha}_f, \tag{17}$$

$$\hat{\alpha}_f = \alpha_r \sqrt{\frac{2}{\alpha_\rho^2 + 1}}, \tag{18}$$

$$\hat{\alpha}_e = 1. \tag{19}$$

This is the basis of our claim that Unit Scaling is what enables a more intuitive set of multipliers. Not only do the multipliers α_r and α_ρ represent important dynamics in the network at initialization (the ratio of residual-to-embedding scales, and the ratio of attention-to-feed-forward scales), but it's only via unit scaling that these equations become simple enough to implement in practice. Using equations Equations (17) to (19) for a non-unit scaled network may still be effective, but the interpretation we've given to α_r and α_ρ no longer hold.

Our final desired property is an empirical one: that the most effective choice of one multiplier depends as little as possible on the choice of the other multiplier(s). We demonstrate that our multipliers satisfy this property better than the standard set of residual multipliers in Section 4.1.

G.2.2 THE FULL U- μ P RESIDUAL SCHEME

Here we give the full definition of our u- μ P residual scheme, summarizing the results of previous sections. A general pre-norm transformer is implemented as:

$$R_0(x) = cx, \quad (20)$$

$$R_l(x) = a_l f_l(R_{l-1}(x)) + b_l R_{l-1}(x), \quad l = 1, \dots, L \quad (21)$$

$$R_{L+1}(x) = f_{L+1}(R_L(x)), \quad (22)$$

where a_l, b_l and c are scalar multipliers, and the f_l alternate between self-attention and feed-forward layers. We consider our baseline set of μ P residual HPs here to be $(\alpha_{\text{emb}}, \alpha_{\text{attn-residual}}, \alpha_{\text{ffn-residual}})$, which we implement (assuming depth- μ P branch scaling) as:

$$a_l = \begin{cases} \frac{\alpha_{\text{attn-residual}}}{\sqrt{L/2}} & l \text{ is odd (self-attention)} \\ \frac{\alpha_{\text{ffn-residual}}}{\sqrt{L/2}} & l \text{ is even (feed-forward)} \end{cases}$$

$$b_l = 1$$

$$c = \alpha_{\text{emb}}.$$

The corresponding u- μ P set of residual HPs is $(\alpha_{\text{res}}, \alpha_{\text{res-attn-ratio}})$, which we implement as:

$$a_l^2 = \frac{\tau_l^2}{\tau_l^2 + 1} \quad (23)$$

$$b_l^2 = \frac{1}{\tau_l^2 + 1} \quad (24)$$

$$c = 1, \quad (25)$$

$$\tau_l^2 = \begin{cases} \frac{\hat{\alpha}_a^2}{\frac{L}{2} + \ell \hat{\alpha}_a^2 + \ell \hat{\alpha}_f^2} & l \text{ is odd} \\ \frac{\hat{\alpha}_f^2}{\frac{L}{2} + (\ell + 1) \hat{\alpha}_a^2 + \ell \hat{\alpha}_f^2} & l \text{ is even} \end{cases}, \quad \ell = \left\lfloor \frac{l-1}{2} \right\rfloor \quad (26)$$

$$\hat{\alpha}_a^2 = \alpha_{\text{res-attn-ratio}}^2 \hat{\alpha}_f^2 \quad (27)$$

$$\hat{\alpha}_f^2 = \frac{2}{\alpha_{\text{res-attn-ratio}}^2 + 1} \alpha_{\text{res}}^2. \quad (28)$$

$$\hat{\alpha}_f^2 = \frac{2}{\alpha_{\text{res-attn-ratio}}^2 + 1} \alpha_{\text{res}}^2. \quad (29)$$

This is the u- μ P residual scheme. It satisfies the three properties that we initially set out to achieve: the variance at initialization of our $R_l(x)$ is always 1, our HPs have a clear and useful interpretation, and our scheme is as expressive as the baseline (which is neither unit-scaled or has interpretable HPs).

H THE CUT-EDGE RULE

In the section we review the notion of *constraints* used for scaling operations in a computational graph. For a more thorough, generalized treatment, please refer to Section 5.1 and Appendix E.4 of the Unit Scaling paper (Blake et al., 2023).

For simplicity, we will only discuss the cut-edge rule in the context of a typical neural network. For each operation f , parametrized by θ taking input x and emitting output y , a user must choose how to scale y , ∇_x and ∇_θ (gradient of loss w.r.t. x and θ respectively). In the simplest case, where there are no further data dependencies, we can simply choose factors that preserve unit scale. In more complex scenarios, we must balance the need for each tensor to be unit-scaled and for gradients to be correct up to a constant factor.

In particular, a problem emerges in the presence of residual blocks in which $y = x + f(x; \theta)$. In these circumstances, ∇_x is computed as the sum of residual gradient $\nabla_f \frac{\partial f}{\partial x}$ and skip gradient ∇_y . If we choose not to insert scaling factors into our graph, $\nabla_f \frac{\partial f}{\partial x}$ and ∇_y will have some ratio of scale r . However, if we have chosen to rescale the gradient of operations in f , then $\nabla_f \frac{\partial f}{\partial x}$ will have been rescaled by some s . This means the new ratio of $\nabla_f \frac{\partial f}{\partial x}$ and ∇_y will be $r \cdot s$. Therefore, when adding these together, ∇_x is no longer a correct gradient up to a constant factor.

How do you remedy this? If we can ensure that the scaling factors are the same for both the input gradients and outputs of an op, we will have $s = 1$. This ensures that gradients for inputs to residual blocks are correct up to a constant factor.

How do you decide when you are free to preserve unit scale, and when to constrain scaling factors to be the same? We previously define the *cut-edge rule* (Blake et al., 2023) for computational graphs where nodes represent forward pass operations and edges represent operation outputs. If an input edge is a *cut-edge*, i.e., the number of connected components in the graph would increase upon deletion (examples in a typical transformer model: output of embedding gather, output of a residual add, output of final norm, output token logits, weights), there is no need to constrain the scales of the operation’s output edge and the input edge gradient. For all other input edges (e.g., inputs to a residual add, intermediates computed along a residual branch), the scales of gradients and outputs should be constrained.

I FROM μ P TO U- μ P

Here we outline additional details to help readers follow the process of deriving u- μ P from the combination of Unit Scaling and μ P. Our first step of dropping σ_W and base-fan-in, and moving α_W s to functions, results in Table 11. This intermediate scheme does not yet satisfy Unit Scaling, but simplifies the HP rules in preparation for further changes. Note that we have also removed $\hat{\eta}_{\text{emb}}$ as we don’t include this HP in our u- μ P extended HP set. We have included residual scaling rules here, in accordance with depth- μ P, which we intend u- μ P to satisfy, though our standard μ P implementation doesn’t use it.

Table 11: An intermediate scheme resulting from dropping those HPs from μ P which are not needed under u- μ P.

ABC-multiplier		Weight Type			Residual
		Input	Hidden	Output	
parameter	(A_W)	1	1	$\frac{1}{\text{fan-in}}$	$\frac{1}{\sqrt{\text{depth}}}$ *
initialization	(B_W)	1	$\frac{1}{\sqrt{\text{fan-in}}}$	1	—
Adam LR	(C_W)	η	$\eta \frac{1}{\text{fan-in}}$	η	$\frac{1}{\sqrt{\text{depth}}}$

J LOW-PRECISION AND ITS TRADE-OFFS

Number formats for deep learning The standard numerical representations used in deep learning are the set of formats defined by the IEEE 754 floating-point standard (IEEE Computer Society, 2019). IEEE floats comprise three elements: a sign bit, exponent bits, and mantissa bits. The number of exponent bits determines the *range* of a format, while the mantissa determines the *precision*⁵.

We refer readers to Blake et al. (2023), Section 3.1 for a comprehensive overview of floating-point representations.

The default format used for training is the single-precision floating-point format, commonly known as FP32, with some hardware providers automatically casting it to the smaller TF32 compute mode for accelerated arithmetic. The 16-bit FP16 and BF16 formats were later introduced, and more recently the FP8 E5 & E4 formats (Sun et al., 2019; Noune et al., 2022; Micikevicius et al., 2022). The higher range of E5 has typically been used for gradients, while the higher precision of E4 has been seen as necessary for weights and activations. Other aspects of training such as the optimizer state and cross-device communication have also been put into FP8 (Peng et al., 2023), though not all tensors are amenable to being run in the lowest precision (Dettmers et al., 2022) without degradation. The use of multiple formats is known as *mixed precision* (Micikevicius et al., 2018). A comparison of these formats is given in Table 12.

Table 12: A comparison of deep learning formats. E indicates exponent bits, and M mantissa bits. The smaller formats typically give more FLOPS, at the expense of reduced range and/or precision.

Format	E	M	max	min normal	min subnormal	FLOPS (vs TF32)
FP32	8	23	3.4×10^{38}	1.2×10^{-38}	1.4×10^{-45}	$< 1 \times$
TF32	8	10	3.4×10^{38}	1.2×10^{-38}	1.1×10^{-41}	$1 \times$
BF16	8	7	3.4×10^{38}	1.2×10^{-38}	9.2×10^{-41}	$2 \times$
FP16	5	10	65504	6.1×10^{-5}	6.0×10^{-8}	$2 \times$
FP8 E5	5	2	57344	6.1×10^{-5}	1.5×10^{-5}	$4 \times$
FP8 E4	4	3	448	1.6×10^{-2}	2.0×10^{-3}	$4 \times$

The benefits of low-precision Using numerical representations with fewer bits facilitates the design of more efficient arithmetic in hardware, typically leading to a linear increase in peak FLOPS (as shown in Table 12). As large-scale training efforts are typically compute-bound due to the size of matmuls (Narayanan et al., 2021), putting the inputs to these operations in low-precision formats has a substantial impact on training efficiency. Low-precision formats also reduce the other two common performance constraints: for memory-bandwidth-bound models they require fewer bits to be transmitted, and for memory-size-bound models they require fewer bits to be stored.

The challenges of low-precision Unfortunately, moving to low-precision formats also increases *quantization error*. For values within the representable range this takes the form of *rounding error*, and for values outside it, *clipping error* (both overflow and underflow). Rounding error tends to be an intrinsic problem: the number of mantissa bits dictates the expected accuracy of representations and this cannot easily be changed. In contrast, clipping error is often eliminated by scaling a tensor so that its values lie within the range of a format. Note that a multiplicative change in values of this kind doesn’t affect the (relative) rounding error, due to the exponential spacing of values. Most research into making low-precision work has focused on the problem of scaling tensors in this way.

Simply casting all tensors to FP16 or FP8 tends to impair training, largely due to clipping error. For FP16, this primarily affects gradients. Micikevicius et al. (2018) address this by introducing a fixed global *loss-scale* HP, which multiplies the loss value in the backward pass, artificially up-scaling gradients to lie within FP16 range. *Automatic loss scaling* (Kuchaiev et al., 2018) builds upon this idea, making the loss-scale a dynamic value that is tuned during training.

The later BF16 format has the same range as FP32, making loss scaling unnecessary. For FP8 no such range-equivalent format can exist, so the problem of clipping error must be addressed. Most FP8 implementations have done so by moving from a global loss-scale to a local scale for each FP8 tensor. In pseudo-code, this takes the form:

```
a = scale(A)
b = scale(B)
A = to_fp8(A / a)
```

⁵ Confusingly, the term *low-precision* tends to indicate using <32 bit-width formats, so in this context *precision* also reflects the number of exponent bits as well as the usual mantissa bits.

```

B = to_fp8(B / b)
C = (a * b) * matmul(A, B)

```

where we assume that `matmul` takes inputs in FP8 and directly produces the output in higher precision.

The result of the `scale()` operation can either be a fixed scale determined before training (Noune et al., 2022), or in the case of Transformer Engine (NVIDIA, 2024a), computed dynamically as a function of the ‘absmax’ of the input tensor (though they introduce a delay across time-steps, to facilitate an efficient fused kernel). Increasing granularity and computing scales dynamically using this kind of method inevitably adds complexity (from both a logical and implementation perspective), as well the potential for computational overhead. Unit Scaling generally avoids the need for `matmul` input scaling.

K BENCHMARKING SCALED MATRIX MULTIPLICATION IMPLEMENTATION IN PYTORCH

Given that the end-goal of leveraging u-mup’s low-precision properties is to speed up training and reduce memory usage, it’s reasonable to ask why we don’t investigate this experimentally. The answer relates to the relative immaturity of the FP8 training software stack - a lack of open, efficient FP8 kernels for compute and communication mean significant additional engineering effort is required to attain expected speedups over the full model.

Here we show that u-μP’s static scaling factors add no overhead to `matmuls` in FP8, and hence ought to be able to reach close to the maximal FP8 throughput attainable for the full model.

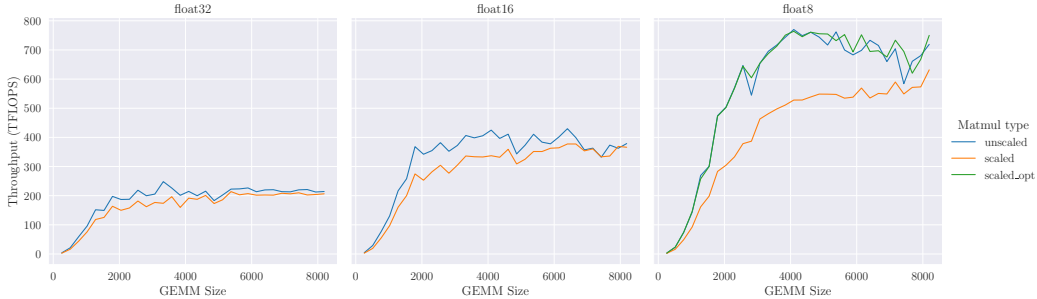


Figure 22: Square matrix multiplication throughput in TFLOPs with and without scaling factors applied to the output across 32-, 16-, and 8-bit float dtypes on NVIDIA H100 PCIe. Naive implementation in PyTorch.

Standard strategies for FP8 training require expensive statistics gathering (e.g., `amax`) per tensor. A key benefit of u-μP for FP8 training is that it instead provides us with static scaling factors to rescale operation outputs. Even a naive implementation in `pytorch` can achieve a minimal drop in hardware utilization.

Figure 22 demonstrates hardware utilization for FP8, FP16, and FP32 matrix multiplications on a single NVIDIA H100 PCIe card. For FP16 and FP32, `torch.matmul` is used, whereas `torch._scaled_mm` is used for FP8. Comparing “scaled” to “unscaled” matrix multiplication demonstrates a 30%, 20%, and 10% drop in hardware utilization for each data type respectively. In the case of FP8, where the drop in utilization is most pronounced, utilization can be recovered by passing the scaling factor as a scale associated with one of the two input tensors.

It should be noted that as of PyTorch version 2.3, `torch._scaled_mm` always computes `amax` as well as the matrix multiplication. The performance of FP8 matrix multiplications could be higher without this overhead.

The above analysis focuses on throughput; significant memory savings are also possible through the use of FP8, though how this affects the total memory footprint depends on various additional variables and the overall distributed training setup. The following factors play a significant role: typically the main memory bottlenecks are the optimizer states, which are kept in full precision. This footprint can be reduced by applying ZeRO sharding (Rajbhandari et al., 2020), though for significant gains the number of data parallel processes needs to be sufficiently large and ZeRO stage 2 or 3 are required. In these settings the memory footprint of activations and gradients becomes significant, and quantizing these to lower precision promises further memory savings, though may be non-trivial (Peng et al., 2023).

L ATTENTION OUTPUT RMS GROWS WITH MODEL DEPTH

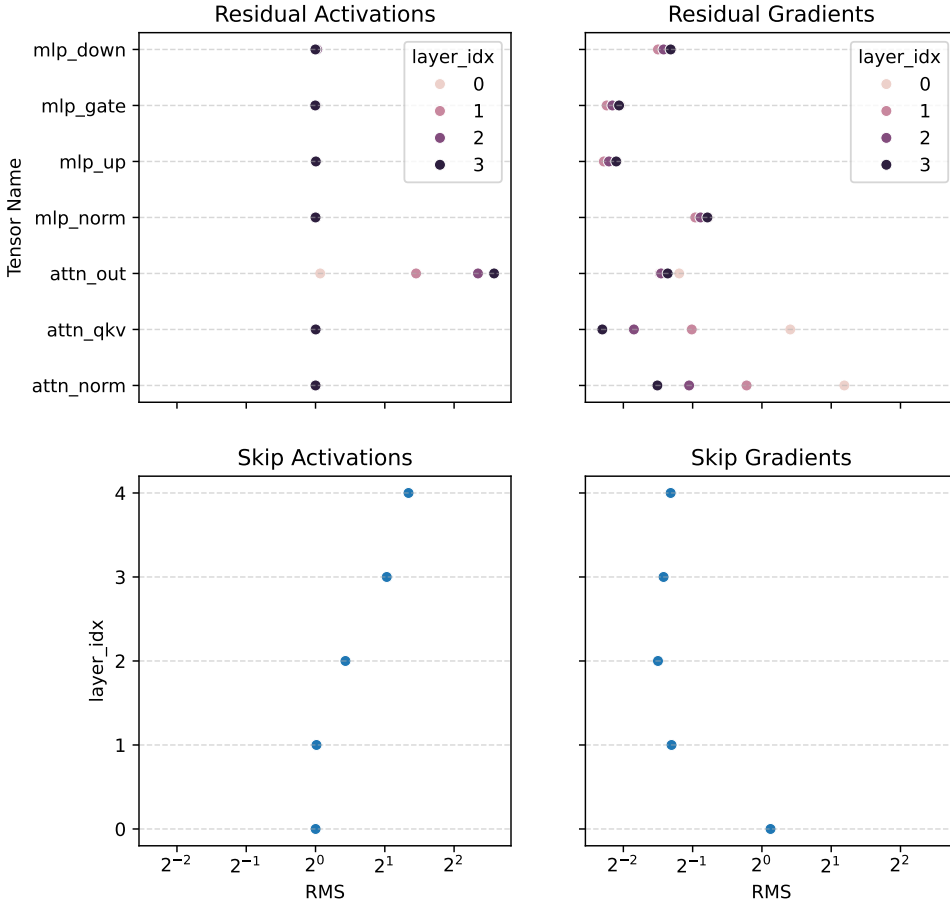


Figure 23: Scale of intermediate tensors grows with depth at initialization. Top left: Intermediate activation tensor RMS along the residual branch. Only the attention outputs after the first layer are not unit-scaled. Bottom left: Skip activation tensor RMS. Scale growth in attention outputs drives growth in skip activation scales. Note that $\text{layer_idx}=0$ corresponds to the embedding output, and $\text{layer_idx}=4$ corresponds to the final layer outputs. Top right: Intermediate gradient tensor RMS along the residual branch. Growth in the attention output scale drives growth in attention qkv gradient scales. Bottom Right: Skip gradient tensor RMS. The scale of output activations induces a global rescaling of the gradients.

A core assumption in deriving per-op scaling factors is that each input to an operation has zero mean, unit-variance, and uncorrelated elements at initialization. This is trivially true for weights and by extension the token embeddings taken as input to the transformer trunk. However, this is not guaranteed for intermediate results and gradients if an operation in the computational graph induces correlation in the elements. In such a scenario our scaling factors will not return unit-variance outputs as we will not have corrected for these correlations in the inputs. As we then increase the depth of the network, where the same operation is left to amplify correlations, we can end up with variance in intermediate results and gradients scaling with depth

Figure 23 illustrates this phenomenon in a unit-scaled four-layer Llama model with width=256. All activation tensors in the residual branches are unit-scaled, except for the output of the attention layers. We also see that the variance of attention outputs grows with depth. Since Llama models use pre-norm on the residual-branch, residual-branch inputs will revert to unit-scale again until they reach another instance of the correlation-inducing operation. As we add under-scaled attention layer results back to the skip-branch, our skip tensor variances grow with depth as our residual-add assumes unit-variance inputs. This has a knock-on effect on the global scaling of the gradients since the Jacobian of the final norm will scale the gradient by the inverse of the final skip tensor variance.

So which operation induces correlation in the attention output at initialization? For the default case where all multipliers are set to 1, our $1/d$ scaling of attention logits results in a sufficiently high temperature that attention probabilities are effectively uniform. With causal masking, we effectively take a running mean across the value tensor along the sequence dimension. As a result, each subsequent token representation is correlated with the last. Since we derive appropriate scaling factors for the first layer, we do not see scale growth emerging until the second layer, where correlations accumulate during the next effective running mean.

We leave it to future work to offer a solution to scale growth created by correlation in intermediate tensors. We note that this is scale growth emergent at initialization, but we also see scale growth in other intermediate tensors during training. Whether scale growth during training is related to the phenomenon outlined here remains to be seen.