

Tuning Large Neural Networks via Zero-Shot Hyperparameter Transfer

Anonymous Author(s)
Affiliation
Address
email

Abstract

Hyperparameter (HP) tuning in deep learning is an expensive process, prohibitively so for neural networks (NNs) with billions of parameters. We show that, in the recently discovered Maximal Update Parametrization (μ P), many optimal HPs remain stable even as model size changes. This leads to a new HP tuning paradigm: parametrize the target model in μ P, tune the HP indirectly on a smaller model, and *zero-shot transfer* them to the full-sized model, i.e., without directly tuning the latter at all. We verify our approach on Transformer and ResNet. For example, by transferring pretraining HPs, we outperform BERT-large on MNLI and QQP, with a total tuning cost (in FLOPs) equivalent to pretraining BERT-large once.

1 Introduction

Hyperparameter tuning is critical to deep learning. Poorly chosen hyperparameters result in subpar performance and training instability. Many published baselines are hard to compare to one another due to varying degrees of hyperparameter tuning. These issues are exacerbated when training extremely large deep learning models, since state-of-the-art networks with billions of parameters become prohibitively expensive to tune.

Recently, [40] showed that different neural network parametrizations induce different infinite-width limits and proposed the *Maximal Update Parametrization (abbreviated μ P)* (summarized in Table 3) that enables “maximal” feature learning in the limit. Intuitively, it ensures that each layer is updated on the same order during training *regardless of width*.¹ We leverage this new parametrization to *zero-shot transfer hyperparameters*

from small models to large models in this work – that is, we obtain near optimal hyperparameters on a large model without directly tuning it at all! While practitioners have always guessed hyperparameters of large models from those of small models, the results are hit-or-miss at best, and this is because of incorrect parametrization. For example, as shown in Fig. 1, in a transformer, the optimal learning rate is stable with width in μ P (right) but far from stable in standard parametrization (left). In addition to width, we empirically verify that, with a few caveats, hyperparameters can also be transferred across depth (in the main text) as well as batch size, language model sequence length, and training time (in the appendix). This reduces the tuning problem of an (arbitrarily) large model to that of a (fixed-sized)

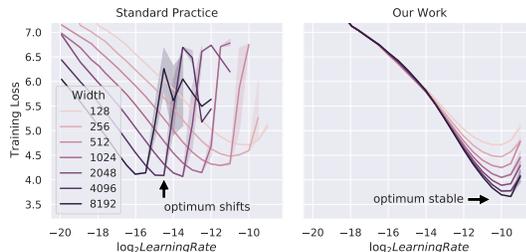


Figure 1: Training loss against learning rate on Transformers of varying d_{model} trained with Adam. Conventionally and in contrast with our technique, different widths do not share the same optimal hyperparameter; wider networks do not always perform better than narrower ones; in fact they underperform the same-width networks in our technique even after tuning learning rate. See Sections 3 and 4 for experimental setup.

¹i.e., updates’ effect on activations become roughly independent of width in the large width limit

Algorithm 1 Tuning a Large Target Model by μ Transfer

- 1: Parametrize target model in Maximal Update Parametrization (μ P)
 - 2: Tune a smaller version (in width and/or depth) of target model
 - 3: Copy tuned hyperparameters to target model
-

Table 1: **Hyperparameters That Can Be μ Transferred, Not μ Transferred, or μ Transferred Across**, with a few caveats discussed in Section 5.1. * means *empirically validated only on Transformers*, while all others additionally have theoretical justification.

| μ Transferable | Not μ Transferable | μ Transferred Across |
|--|---|---|
| optimization related, init, parameter multipliers, etc | regularization (dropout, weight decay, etc) | width, depth*, batch size*, training time*, seq length* |

37 small model. Our overall procedure, which we call μ Transfer, is summarized in Algorithm 1 and
38 Fig. 2, and the hyperparameters we cover are summarized in Tables 1 and 2.

39 There are several benefits to our approach: 1. **Speedup:** It
40 provides massive speedup to the tuning of large models.
41 For example, we are able to outperform published numbers
42 of BERT-large [7], as measured on MNLI and QQP [37],
43 purely by zero-shot hyperparameter transfer, with tuning
44 cost approximately equal to only 1 BERT-large pretraining.
45 For the largest models such as T5 [27] or GPT-3 [6], our
46 approach renders hyperparameter tuning possible at all.²

47 2. **Tune once for whole family:** For any fixed family of
48 models with varying width and depth (such as the BERT
49 family or the GPT-3 family), we only need to tune a single
50 small model and can reuse its hyperparameters for all
51 models in the family (but in general not for different data
52 and/or tasks). For example, we will use this technique to
53 tune BERT-base and BERT-large simultaneously. 3. **Better**

54 **Compute Utilization:** While large model training needs to be distributed across many GPUs, the
55 small model tuning can happen on individual GPUs, greatly increasing the level of parallelization of
56 tuning (and in the context of organizational compute clusters, better scheduling and utilization ratio).

57 Nevertheless, μ Transfer still has several limitations. For example, while it is very effective for
58 pretraining, it cannot transfer regularization hyperparameters, so it’s generally not applicable to the
59 finetuning of pretrained models. We discuss other limitations carefully in Section 5.1.

60 Our Contributions

- 61 • We demonstrate it is possible to zero-shot transfer near optimal hyperparameters to a large
62 model from a small version;
- 63 • We propose a new hyperparameter tuning technique, μ Transfer, for large neural networks
64 based on this observation that provides massive speedup over conventional methods;
- 65 • We thoroughly verify our method on machine translation and large language models pre-
66 training (in main text) as well as image classification (in appendix);
- 67 • We release a Pytorch[24] package for implementing μ Transfer painlessly. A sketch of this
68 package is given in Appendix E.

69 **Terminologies** Sometimes, to be less ambiguous, we will often refer to the “large model” as the
70 *target model*, as it is the model we wish to ultimately tune, while we refer to the “small model” as
71 the *proxy model*, as it proxies the hyperparameter tuning process. We will follow standard notation
72 $d_{model}, d_{head} = d_k, d_v, n_{head}, d_{ffn}$ regarding dimensions in a Transformer; one can see Fig. 6 for a
73 refresher.

²Note, again, that our tuning cost stays *fixed* even as the target model grows in size, so tuning T5-large would have the same cost as tuning BERT-large even though it is 4 times larger.

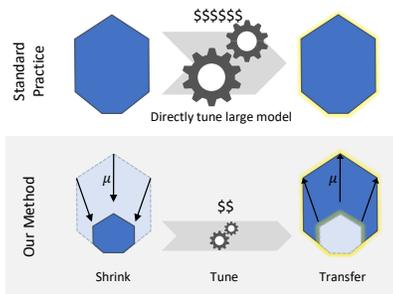


Figure 2: Illustration of μ Transfer

Table 2: **Examples of μ Transferable Hyperparameters.** All of the below can also be specialized to per-layer hyperparameters.

| Optimizer Related | Initialization | Parameter Multipliers |
|---|--------------------------|---|
| learning rate (LR), momentum, Adam beta, LR schedule, etc | per-layer init. variance | multiplicative constants after weight/biases, etc |

74 2 Parametrization Matters: A Primer

75 In this section, we give a very basic primer on why the correct parametrization can allow hyperpa-
 76 rameter transfer across width, but see Appendices G.1 to G.3 for more (mathematical) details.

77 The Central Limit Theorem (CLT) says that, if x_1, \dots, x_n are iid samples from a zero-mean, unit-
 78 variance distribution, then $\frac{1}{\sqrt{n}}(x_1 + \dots + x_n)$ converges to a standard Gaussian $\mathcal{N}(0, 1)$ as $n \rightarrow \infty$.

79 Therefore, we can say that $\frac{1}{\sqrt{n}}$ is the right order of *scaling factor* c_n such that $c_n(x_1 + \dots + x_n)$
 80 converges to something nontrivial. In contrast, if we set $c_n = 1/n$, then $c_n(x_1 + \dots + x_n) \rightarrow 0$; or
 81 if $c_n = 1$, then $c_n(x_1 + \dots + x_n)$ blows up in variance as $n \rightarrow \infty$.

82 Now suppose we would like to minimize the function

$$F_n(c) \stackrel{\text{def}}{=} \mathbb{E}_{x_1, \dots, x_n} f(c(x_1 + \dots + x_n)) \quad (1)$$

83 over $c \in \mathbb{R}$, for some bounded continuous function $f : \mathbb{R} \rightarrow \mathbb{R}$. If we reparametrize $c = \alpha/\sqrt{n}$ for
 84 $\alpha \in \mathbb{R}$, then by CLT, $G_n(\alpha) \stackrel{\text{def}}{=} F_n(c) \rightarrow \mathbb{E} f(\mathcal{N}(0, \alpha^2))$ stabilizes into a function of α as $n \rightarrow \infty$.

85 Then for sufficiently large n , the optimal $\alpha_n^* \stackrel{\text{def}}{=} \arg \min_{\alpha} G_n(\alpha)$ should be close to α_N^* for any
 86 $N > n$, and indeed, for $N = \infty$ — this precisely means we can *transfer* the optimal c_n^* or α_n^* for a
 87 smaller problem (say F_n) to a larger problem (say F_N): G_N is approximately minimized by α_n^* and
 88 F_N is approximately minimized by $c_n^* \sqrt{N/n}$. Because the transfer algorithm is simply copying α ,
 89 we say the parametrization $c = \alpha/\sqrt{n}$ is the *correct parametrization* for this problem.

90 In the scenario studied in this paper, x_1, \dots, x_n are akin to randomly initialized parameters of a
 91 width- n neural network, c is akin to a hyperparameter such as learning rate, and f is the test-set
 92 performance of the network *after training*, so that F_n gives its expectation over random initializations.
 93 Just as in this example, if we parametrize the learning rate and other hyperparameters correctly,
 94 then we can directly copy the optimal hyperparameters for a narrower network into a wide network
 95 and expect approximately optimal performance — this is the *hyperparameter transfer* we propose
 96 here. It turns out the Maximal Update Parametrization (μ P) introduced in [40] is correct (akin to
 97 the parametrization in α above), while the standard parametrization (SP) is incorrect (akin to the
 98 parametrization in c). We will review both parametrizations shortly. Theoretically, a μ P network has
 99 a well-defined infinite-width limit — akin to $(x_1 + \dots + x_n)/\sqrt{n}$ having a $\mathcal{N}(0, 1)$ limit by CLT —
 100 while a SP network does not (the limit will blow up) [40].³ More concretely, as shown in [40] and
 101 Appendix G.3, in SP, the last layer is initialized too large and is updated disproportionately more as the
 102 model width increases, forcing a smaller learning rate to prevent a blow-up, consequently sacrificing
 103 performance.

104 3 Hyperparameters Don’t Transfer Conventionally

105 In the community there seem to be conflicting assumptions about hyperparameter stability. *A priori*,
 106 models of different sizes don’t have any reason to share the optimal hyperparameters. Indeed, papers
 107 aiming for state-of-the-art results often tune them separately. On the other hand, a nontrivial fraction
 108 of papers in deep learning fixes all hyperparameters when comparing against baselines, which reflects
 109 an assumption that the optimal hyperparameters should be stable — not only between the same model
 110 of different sizes but also between models of different designs — so that such comparisons are fair.
 111 Here, we demonstrate hyperparameter *instability* across width explicitly in MLP and Transformers in
 112 the standard parametrization. We will only look at training loss to exclude the effect of regularization.

³The more theoretically astute reader may observe that SP with a $\Theta(1/\text{width})$ learning rate induces a well-defined infinite-width limit exists as well. Nevertheless, this does not allow hyperparameter transfer because this limit is in kernel regime as shown in [40]. See Appendix G.3 for more discussions.

113 **MLP with Standard Parametrization** We start with a 2-hidden-layer MLP with activation function ϕ , using the standard parametrization⁴ with LeCun initialization⁵ akin to the default in PyTorch:
 114
 115

$$f(\xi) = W^{3\top} \phi(W^{2\top} \phi(W^{1\top} \xi + b^1) + b^2) \quad (2)$$

with init. $W^1 \sim \mathcal{N}(0, 1/d_{in})$, $W^{\{2,3\}} \sim \mathcal{N}(0, 1/n)$, $b^{\{1,2\}} = 0$,

116 where $W^1 \in \mathbb{R}^{d_{in} \times n}$, $b^1 \in \mathbb{R}^n$,
 117 $W^2 \in \mathbb{R}^{n \times n}$, $b^2 \in \mathbb{R}^n$, $W^3 \in$
 118 $\mathbb{R}^{n \times d_{out}}$ and d_{in} , n , and d_{out} are
 119 the input, hidden, and output dimensions. The particular MLP we use has
 120 $\phi = ReLU$ and a cross-entropy (xent) loss function. We define the width of
 121 MLP as the hidden size n , which is
 122 varied from 256 to 8192. The models
 123 are trained on CIFAR-10 for 20
 124 epochs, which is more than enough to
 125 ensure convergence.
 126

128 As shown on the left in Fig. 3, the
 129 optimal learning rate shifts by roughly
 130 an order of magnitude as the width
 131 increases from 256 to 8192; using the
 132 optimal learning of the smallest model
 133 on the largest model gives very bad performance, if not divergence.

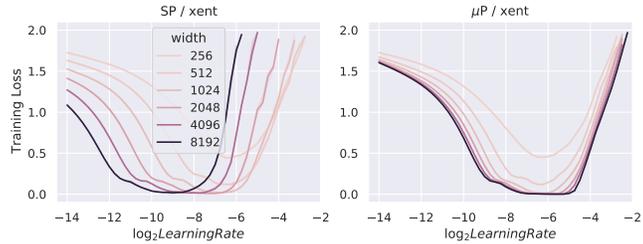


Figure 3: MLP width different hidden sizes trained for 20 epoch on CIFAR-10 using SGD. **Left** uses standard parametrization (SP); **right** uses maximal update parametrization (μ P). μ P networks exhibit better learning rate stability than their SP counterparts.

134 **Transformer with Standard Parametrization** This perhaps unsurprising observation holds for
 135 more complex architectures such as Transformer as well, as shown in Fig. 1 (left). We define width
 136 as d_{model} , with $d_k = d_q = d_v = d_{model}/n_{heads}$ and $d_{ffn} = 4d_{model}$. The models are trained on
 137 wikitext-2 for 5 epochs. In Fig. 12 in the appendix we also show the instability of initialization scale
 138 and other hyperparameters.

139 4 Unlocking Zero-Shot Hyperparameter Transfer with μ P

140 We show that μ P solves the problems we see in Section 3.

141 **MLP with μ P** The *basic form* of μ P for the MLP in Section 3 is

$$f(\xi) = 1/\sqrt{n} W^{3\top} \phi(W^{2\top} \phi(\sqrt{n} W^{1\top} \xi + \sqrt{n} b^1) + \sqrt{n} b^2) \quad (3)$$

with init. $W^1 \sim \mathcal{N}(0, 1/n)$, $W^{\{2,3\}} \sim \mathcal{N}(0, 1/n)$, $b^{\{1,2\}} = 0$.

142 Here we highlighted in purple the differences in the two parametrizations, namely the first layer
 143 initialization and the explicit multipliers in front of first layer weights, biases, and the output in
 144 μ P. This basic form makes clear the *scaling with width* n of the parametrization, but in practice we
 145 will often insert (possibly tune-able) multiplicative constants in front of each appearance of n . For
 146 example, a particularly useful instance of this is when we would like to be consistent with a SP MLP
 147 at a *base width* n_0 . Then we may insert constants as follows: For $\tilde{n} \stackrel{\text{def}}{=} n/n_0$,

$$f(\xi) = 1/\sqrt{\tilde{n}} W^{3\top} \phi(W^{2\top} \phi(\sqrt{\tilde{n}} W^{1\top} \xi + \sqrt{\tilde{n}} b^1) + \sqrt{\tilde{n}} b^2) \quad (4)$$

with init. $W^1 \sim \mathcal{N}(0, 1/\tilde{n} \cdot 1/d_{in})$, $W^{\{2,3\}} \sim \mathcal{N}(0, 1/n)$, $b^{\{1,2\}} = 0$,

148 Then at width $n = n_0$, all purple factors above are 1, and the parametrization is identical to SP
 149 (Eq. (2)) at width n_0 . Of course, as n increases from n_0 , then Eq. (4) quickly deviates away from
 150 Eq. (2). In other words, for a particular n , μ P and SP can be identical up to the choice of some
 151 constants (in this case n_0), but μ P determines a different “set” of networks than SP as one varies n .
 152 As we will see, this deviation is crucial for hyperparameter transfer.

153 Indeed, in Fig. 3(right), we plot the CIFAR10 performances, over various learning rates and widths,
 154 of μ P MLPs with $n_0 = 128$. In contrast to SP, the optimal learning rate under μ P is stable. This

⁴i.e. the default parametrization offered by common deep learning frameworks. See Table 3 for a review.

⁵The key here is that the init. variance $\propto 1/\text{fan_in}$, so the same insights here apply with e.g. He initialization.

Table 3: $\mu\mathbf{P}$ [40] and SP for General Neural Networks, Basic Form. This basic form emphasizes the *scaling with width* (fan_in or fan_out); in practice, we may insert tunable multipliers in front of fan_in and fan_out as in Eq. (4). Notations: 1) η is the “master” learning rate. 2) The fan_out of a bias vector is its dimension (whereas fan_in is 1). 3) *Multiplier* means explicit multipliers of the parameter, such as in Eq. (3). 4) **Purple text** highlights key differences from standard parametrization (SP); **Gray text** recalls the corresponding SP. *SGD* (resp. *Adam*) here can be replaced by variants such as SGD with momentum (resp. Adagrad, Adadelata, etc). Transformer $\mu\mathbf{P}$ requires one more modification ($1/d$ attention instead of $1/\sqrt{d}$); see Definition 4.1.

| | Input weights & all biases | | Output weights | | Hidden weights | |
|------------|-------------------------------|----------------------|------------------------------|----------|-----------------------|----------|
| Init. Var. | $1/\text{fan_out}$ | $(1/\text{fan_in})$ | $1/\text{fan_in}$ | | $1/\text{fan_in}$ | |
| Multiplier | $\sqrt{\text{fan_out}}$ | (1) | $1/\sqrt{\text{fan_in}}$ | (1) | 1 | |
| SGD LR | | η | η | | η | |
| Adam LR | $\eta/\sqrt{\text{fan_out}}$ | (η) | $\eta/\sqrt{\text{fan_in}}$ | (η) | $\eta/\text{fan_in}$ | (η) |

155 means that, the best learning rate for a width-128 network is also best for a width-8192 network in
 156 $\mu\mathbf{P}$ — i.e. hyperparameter transfer *works* — but not for SP. In addition, we observe performance for a
 157 fixed learning rate always increases with width in $\mu\mathbf{P}$, but not in SP.

158 This MLP $\mu\mathbf{P}$ example can be generalized easily to general neural networks trained under SGD or
 159 Adam, as summarized in Table 3.

160 **Transformers with $\mu\mathbf{P}$** We repeat the experiments with base width $n_0 = 128$ for Transformers:

161 **Definition 4.1.** The *Maximal Update Parametrization* ($\mu\mathbf{P}$) for a Transformer is given by Table 3
 162 and $1/d$ attention instead of $1/\sqrt{d}$, i.e. the attention logit is calculated as $q^\top k/d$ instead of $q^\top k/\sqrt{d}$
 163 where query q and key k have dimension d .⁶

164 The results are shown on the right in Fig. 1, where the optimal learning rate is stable, and the
 165 performance improves monotonically as width increases.

166 5 Which Hyperparameters Can Be μ Transferred?

167 In this section, we explore how common hyperparameters fit into our framework. In general, they can
 168 be divided into three kinds, summarized in Table 1:

- 169 1. those that can transfer from the small to the large model, such as learning rate (Table 2);
- 170 2. those that primarily control regularization and don’t work well with our technique; and
- 171 3. those that define training *scale*, such as width as discussed above as well as others like depth
 172 and batch size, across which we transfer other hyperparameters.

173 Those in the first category transfer across width, as theoretically justified above in Section 2, while
 174 we empirically explore the transfer across the other dimensions in the third category, in order to push
 175 the practicality and generality of our technique. Note that μ Transfer across width is quite general,
 176 e.g. it allows varying width ratio of different layers or number of attention heads in a Transformer;
 177 see Appendix B.2. This will be very useful in practice. For the second category, the amount of
 178 regularization naturally depends on both the model size and data size, so we should not expect transfer
 179 to work if the parametrization only depends on model size. We discuss these hyperparameters in
 180 more detail in Appendix B.1.

181 5.1 Empirical Validation and Limitations

182 Our empirical investigations focus on Transformers (here) and ResNet (in Appendix D.1.1), the most
 183 popular backbones of deep learning models today. We train a 2-layer pre-layernorm $\mu\mathbf{P}$ ⁷ Transformer

⁶This is roughly because during training, q and k will be correlated so $q^\top k$ actually scales like d due to Law of Large Numbers, in contrast to the original motivation that q, k are uncorrelated at initialization so Central Limit applies instead. See Appendix G.2.1 for a more in-depth discussion.

⁷“2 layers” means the model has 2 self-attention blocks. To compare with SP Transformer, see Fig. 12.

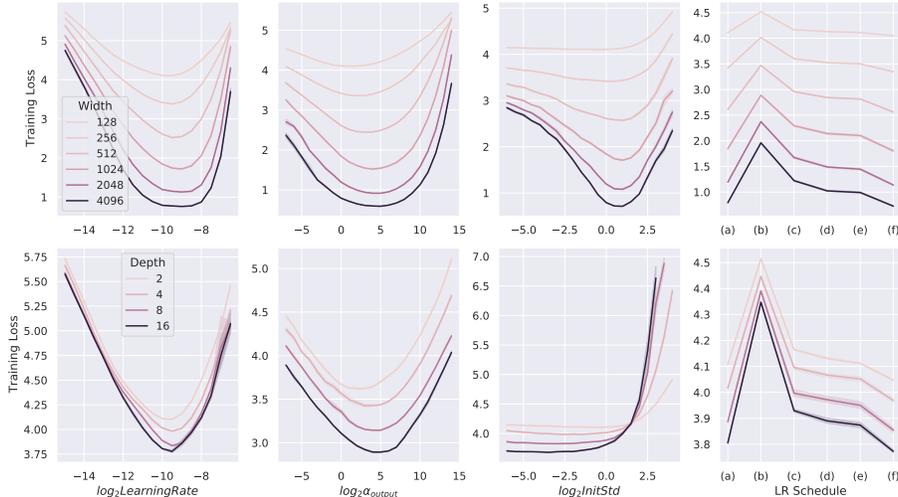


Figure 4: **Empirical validation of the stability of four representative hyperparameters on pre-LN Transformers in μ P**: learning rate, last layer weight multiplier α_{output} , weight initialization standard deviation, and learning rate schedule. We use the following learning rate schedules: (a) constant; (b) linear decay; (c) StepLR @ [5k, 8k] with a decay factor of 0.1; (d) StepLR @ [4k, 7k] with a decay factor of 0.3; (e) cosine annealing; (f) inverse square-root decay. All models are trained on wikitext-2 for 10k steps. When not specified in the legend, the width used is 256, depth 2, batch size 20, sequence length 32, LR schedule constant. We sweep a particular hyperparameter, corresponding to each column, while fixing all others constant. See Section 5.1 for discussion of these results.

184 with 4 attention heads on Wikitext-2. We sweep one of four hyperparameters (learning rate, output
 185 weight multiplier α_{output} , initialization standard deviation, and learning rate schedule) while fixing
 186 the others and sweeping along width and depth (with additional results in Fig. 10 on transfer across
 187 batch size, sequence length, and training time). Fig. 4 shows the results averaged over 5 random
 188 seeds.

189 Empirically, we find that for language modeling on Transformers, hyperparameters generally transfer
 190 across scale dimensions if some minimum width (e.g. 256), depth (e.g., 4), batch size (e.g., 32),
 191 sequence length (e.g., 128), and training steps (e.g., 5000) are met, with some caveats discussed
 192 below. While the exact optimum can shift slightly with increasing scale, this shift usually has very
 193 small impact on the loss, compared to SP (Figs. 1 and 3(left)). However, there are some caveats.
 194 For example, the best initialization standard deviation does not seem to transfer well across depth
 195 (2nd row, 3rd column), despite having a stabler optimum across width. In addition, while our results
 196 on width, batch size, sequence length, and training time still hold for post-layernorm (Fig. 11),⁸
 197 the transfer across depth only works for pre-layernorm Transformer. Nevertheless, in practice (e.g.
 198 our results in Section 6.2.1) we find that fixing initialization standard deviation while tuning other
 199 hyperparameters works well when transferring across depth.

200 6 Efficiency and Performance of μ Transfer

201 Now that the plausibility of μ Transfer has been established in toy settings, we turn to more realistic
 202 scenarios to see if one can achieve tangible gains. Specifically, we perform hyperparameter tuning
 203 only on a smaller proxy model, test the obtained hyperparameters on the large target model directly,
 204 and compare against baselines tuned using the target model. We seek to answer the question: Can
 205 μ Transfer make hyperparameter tuning more efficient while achieving performance on par with
 206 traditional tuning? As we shall see by the end of the section, the answer is positive. We focus on
 207 Transformers here, while experiments on ResNets on CIFAR10 and Imagenet can be found as well in
 208 Appendix D.1. All of our experiments are run on V100 GPUs.

⁸in fact, post-layernorm transformers are much more sensitive to hyperparameters than pre-layernorm, so our technique is more crucial for them, especially for transfer across width. Fig. 1 uses post-layernorm.

Table 4: **Transformer on IWSLT14 De-En.** 1x and 0.25x refers to scaling of width only. Compared to traditional tuning (“Tuning on 1x”), hyperparameter transfer from 0.25x provides better and more reliable outcome given fixed amount of compute. On the other hand, naive transfer (i.e. with SP instead of μ P) fails completely. The percentiles are over independent trials, with each trial involving the entire tuning process with a new hyperparameter random search.

| Setup | Total Compute | #Samples | Val. BLEU Percentiles | | | |
|----------------------------------|---------------|----------|-----------------------|--------------|--------------|--------------|
| | | | 25 | 50 | 75 | 100 |
| fairseq[22] default | - | - | - | - | - | 35.40 |
| Tuning on 1x | 1x | 5 | 33.62 | 35.00 | 35.35 | 35.45 |
| Naive transfer from 0.25x | 1x | 64 | training diverged | | | |
| μ Transfer from 0.25x (Ours) | 1x | 64 | 35.27 | 35.33 | 35.45 | 35.53 |

209 6.1 Transformer on IWSLT14 De-En

210 **Setup** IWSLT14 De-En is a well-known machine translation benchmark. We use the default IWSLT
 211 (post-layernorm) Transformer implemented in fairseq [22] with 40M parameters, which we denote
 212 as the *1x model*.⁹ For μ Transfer, we tune on a *0.25x model* with 1/4 of the width, amounting to 4M
 213 parameters. For this experiment, we tune via random search the learning rate η , the output layer
 214 parameter multiplier α_{output} , and the attention key-projection weight multiplier α_{attn} . See the grid
 215 and other experimental details in Appendix C.1.

216 We compare transferring from the
 217 0.25x model with tuning the 1x model
 218 while controlling the total tuning bud-
 219 get in FLOPs.¹⁰ To improve the repro-
 220 ducibility of our result: 1) we repeat
 221 the entire hyperparameter search pro-
 222 cess (*a trial*) 25 times for each setup,
 223 with number of samples as indicated
 224 in Table 4, and report the 25th, 50th,
 225 75th, and 100th percentiles; 2) we
 226 evaluate each selected hyperparamete-
 227 r combination using 5 random initial-
 228 izations and report the mean perfor-
 229 mance.¹¹

230 We pick the hyperparameter combi-
 231 nation that achieves the lowest valida-
 232 tion loss¹² for each trial. The re-
 233 ported best outcome is chosen accord-
 234 ing to the validation loss during tun-
 235 ing. We compare against the default in
 236 fairseq, which is presumably heav-
 237 ily tuned. The result is shown in Ta-
 238 ble 4.

239 **Performance Pareto Frontier** The result above only describes a particular compute budget. Is
 240 μ Transfer still preferable when we have a lot more (or less) compute? To answer this question, we
 241 produce the compute-performance Pareto frontier in Fig. 5(left), where we repeat the above experi-
 242 ment with different compute budgets. Evidently, our approach completely dominates conventional
 243 tuning.

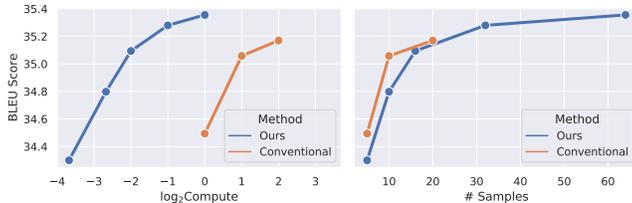


Figure 5: **Efficiency-Performance Pareto frontier** of μ Transfer compared to conventional tuning, on IWSLT Transformer, using random hyperparameter search as the base method. We plot the *median* BLEU score over 25 trials (Left) against relative compute budget in log scale and (Right) against number of hyperparameter samples taken. While with the same number of samples, μ Transfer slightly underperforms conventional tuning, this gap vanishes with more samples, and in terms of compute, our Pareto frontier strongly and consistently dominates that of conventional tuning. Note that, in larger models (e.g. BERT or GPT-3, not shown here), we believe our efficiency advantage will only widen as our small proxy model can stay the same size while the target model grows.

⁹<https://github.com/pytorch/fairseq/blob/master/examples/translation/README.md>.

¹⁰Ideally we would like to measure the wall clock time of tuning. But for smaller models such as the IWSLT Transformer proxy models here, CUDA is poorly optimized, so wall clock time scaling would not reflect the scaling for larger models like BERT. Thus, we measure in FLOPs, which is less dependent on model size.

¹¹We do not report the standard deviation over random initializations to avoid confusion.

¹²We find this provides more reliable result than selecting for the best BLEU score because loss is “smoother.”

244 **Sample Quality of Proxy Model vs Target Model** The Pareto frontier in Fig. 5(right) suggests
 245 that, given a fixed number of random *samples* from the hyperparameter space, 1) tuning the target
 246 model directly yields slightly better results than tuning the proxy model (while taking much more
 247 compute of course), but 2) this performance gap seems to vanish as more samples are taken. This can
 248 be explained by the intuition that the narrower proxy model is a “noisy estimator” of the wide target
 249 model [40].¹³ With few samples, this noise can distort the random hyperparameter search, but with
 250 more samples, this noise is suppressed.¹⁴

251 6.2 Transformer on WMT14 En-De

252 We scale up to WMT14 En-De using the large (post-layernorm) Transformer from [35] with 211M
 253 parameters. We tune on a proxy model with 15M parameters by shrinking d_{model} , d_{fn} , and n_{heads} .
 254 For this experiment, we tune via random search the learning rate η , the output layer parameter
 255 multiplier α_{output} , and the attention key-projection weight multiplier α_{attn} following the grid in
 256 Appendix C.2. The result is shown in Table 5: While random search with 3 hyperparameter samples
 257 far underperforms the fairseq default, we are able to match it via transfer using the same tuning
 258 budget.

Table 5: **Transformers on WMT14 En-De.** 1x and 0.25x refers to scaling of width only. We report BLEU fluctuation over 3 independent trials, i.e., 3 independent random hyperparameter searches.

| Setup | Total Compute | #Samples | Val. BLEU Percentiles | | |
|----------------------------------|---------------|----------|-----------------------|--------------|--------------|
| | | | Worst | Median | Best |
| fairseq[22] default | - | - | - | - | 26.40 |
| Tuning on 1x | 1x | 3 | training diverged | | 25.69 |
| Naive transfer from 0.25x | 1x | 64 | training diverged | | |
| μ Transfer from 0.25x (Ours) | 1x | 64 | 25.94 | 26.34 | 26.42 |

259 6.2.1 BERT

260 Finally, we consider large-scale language model pretraining where hyperparameter tuning is known
 261 to be challenging. Using Megatron (pre-layernorm) BERT [30] as a baseline, we hope to recover the
 262 performance of the published hyperparameters by only tuning a proxy model that has roughly 13M
 263 parameters, which we call *BERT-prototype*. While previous experiments scaled only width, here we
 264 will also scale depth, as discussed in Section 5 and validated in Fig. 4. We use a batch size of 256 for
 265 all runs and follow the standard finetuning procedures. For more details on BERT-prototype, what
 266 hyperparameters we tune, and how we finetune the trained models, see Appendix C.3.

267 During hyperparameter tuning, we sample 256 combinations from the search space and train each
 268 combination on BERT-prototype for 10^5 steps. The total tuning cost measured in FLOPs is roughly the
 269 same as training 1 BERT-large for the full 10^6 steps; the exact calculation is shown in Appendix C.3.
 270 The results are shown in Table 6. Notice that on BERT-large, we obtain sizeable improvement over
 271 the pretuned Megatron BERT-large baseline.

272 7 Related Works

273 **Hyperparameter Tuning** Many have sought to speedup hyperparameter tuning beyond the simple
 274 grid or random search, such as via Bayesian optimization [32, 33] or multi-arm bandits [13, 16].
 275 There are also dedicated tools such as Optuna [4] and Talos [3] which integrate with existing deep
 276 learning frameworks and provide an easy way to apply more advanced tuning techniques. Our work
 277 is complementary to the above, as they can be used to tune the proxy model. it is only for scientific
 278 reasons that we primarily did random search throughout this work.

279 **Hyperparameter Transfer** Many previous works explored transfer learning of hyperparameter
 280 tuning (e.g. [10, 25, 34, 43]). However, to the best of our knowledge, our work is the first to explore

¹³More precisely, the proxy and the target models are both “estimators” of their common *infinite-width limit*, but the former is more noisy than the latter, with width akin to the number of items in an average.

¹⁴but perhaps not completely, as the narrow proxy models may be biased estimators, i.e. the optimal hyperparameters might differ slightly from the wide model as seen in Fig. 4.

Table 6: **BERT pretraining.** Hyperparameter transfer outperforms published baselines without tuning the full model directly at all. We tune BERT-base and BERT-large simultaneously via a single proxy model, *BERT-prototype*. The total tuning cost = the cost of pretraining a single BERT-large. *Model speedup* refers to the training speedup of BERT-prototype over BERT-base or BERT-large. *Total speedup* in addition includes time saving from transferring across training steps. Both speedups can be interpreted either as real-time speedup on V100s or as FLOPs speedup (which turn out to be empirically very similar in this case).

| Model | Method | Model Speedup | Total Speedup | Test loss | MNLI (m/mm) | QQP |
|-----------------------|-----------------------|---------------|---------------|--------------|-------------------|-------------|
| BERT _{base} | Megatron Default | 1x | 1x | 1.995 | 84.2/84.2 | 90.6 |
| BERT _{base} | Naive Transfer | 4x | 40x | | training diverged | |
| BERT _{base} | μ Transfer (Ours) | 4x | 40x | 1.970 | 84.3/84.8 | 90.8 |
| BERT _{large} | Megatron Default | 1x | 1x | 1.731 | 86.3/86.2 | 90.9 |
| BERT _{large} | Naive Transfer | 22x | 220x | | training diverged | |
| BERT _{large} | μ Transfer (Ours) | 22x | 220x | 1.683 | 87.0/86.5 | 91.4 |

281 *zero-shot* hyperparameter transfer. In addition, we focus on transferring across model scale rather
 282 than between different tasks or datasets. Some algorithms like Hyperband [17] can leverage cheap
 283 estimates of hyperparameter evaluations (like using a small model to proxy a large model) but they
 284 are not zero-shot algorithms, so would still be very expensive to apply to large model training.
 285 Nevertheless, all of the above methods are complementary to ours as they can be applied to the tuning
 286 of our proxy model.

287 **Previously Proposed Scaling Rules of Hyperparameters** [9, 21, 29, 31] investigated the right
 288 way to scale learning rate with batch size, with sometimes conflicting proposals. which we summarize
 289 in Appendix A. [23] studied how learning rate (and batch size) should scale with width for MLPs and
 290 CNNs trained with SGD in NTK or standard parametrizations. We provide a detailed comparison of
 291 our work with theirs in Appendix A.

292 Many previous works proposed different initialization or parametrizations with favorable properties,
 293 such as better stability for training deep neural networks [5, 8, 11, 19, 28, 41, 42, 45]. Our work
 294 differs from these in that we focus on the transferability of optimal hyperparameters from small
 295 models to large models in the same parametrization.

296 8 Conclusion

297 Leveraging the discovery of a feature learning neural network infinite-width limit, we hypothesized
 298 and verified that the hyperparameter landscape across NNs of different width is reasonably stable if
 299 parametrized according to Maximal Update Parametrization (μ P). We further empirically showed
 300 that it’s possible to transfer across depth, batch size, sequence length, and training time, with a few
 301 caveats. This allowed us to indirectly tune a very large network by tuning its smaller counterparts
 302 and transferring the hyperparameters to the full model.

303 Nevertheless, our method has plenty of room to improve. For example, initialization does not transfer
 304 well across depth, and depth transfer generally still does not work for post-layernorm Transformers.
 305 This begs the question whether a more principled parametrization in depth could solve these problems.
 306 Additionally, Fig. 4 shows that the optimal hyperparameter still shifts slightly for smaller models.
 307 Perhaps by considering finite-width corrections to μ P one can fix this shift. Finally, it will be
 308 interesting to consider if there’s a way to transfer regularization hyperparameters as a function of
 309 both the model size and data size.

310 **Broader Impact** Our work makes hyperparameter tuning of large models more efficient. This
 311 enables large models to be better tuned given the same compute budget, thereby increasing the
 312 performance per cost. Organizations large and small can focus their research on small models and
 313 scale up only once with reasonable confidence that the training would go well. We do not foresee any
 314 direct negative societal impact.

315 **References**

- 316 [1] NVIDIA/DeepLearningExamples, apache v2 license. URL [https://github.com/NVIDIA/](https://github.com/NVIDIA/DeepLearningExamples)
317 [DeepLearningExamples](https://github.com/NVIDIA/DeepLearningExamples).
- 318 [2] Davidnet, mit license, 2019. URL <https://github.com/davidcpage/cifar10-fast>.
- 319 [3] Autonomio talos, mit license, 2019. URL <http://github.com/autonomio/talos>.
- 320 [4] Takuya Akiba, Shotaro Sano, Toshihiko Yanase, Takeru Ohta, and Masanori Koyama. Optuna:
321 A next-generation hyperparameter optimization framework, 2019.
- 322 [5] Thomas Bachlechner, Bodhisattwa Prasad Majumder, Huanru Henry Mao, Garrison W. Cot-
323 trell, and Julian McAuley. ReZero is All You Need: Fast Convergence at Large Depth.
324 *arXiv:2003.04887 [cs, stat]*, June 2020. URL <http://arxiv.org/abs/2003.04887>. arXiv:
325 2003.04887.
- 326 [6] Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal,
327 Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel
328 Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M.
329 Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz
330 Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec
331 Radford, Ilya Sutskever, and Dario Amodei. Language models are few-shot learners, 2020.
- 332 [7] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: Pre-training of
333 Deep Bidirectional Transformers for Language Understanding. *arXiv:1810.04805 [cs]*, May
334 2019. URL <http://arxiv.org/abs/1810.04805>. arXiv: 1810.04805 version: 2.
- 335 [8] Xavier Glorot and Yoshua Bengio. Understanding the difficulty of training deep feedforward
336 neural networks. In Yee Whye Teh and Mike Titterton, editors, *Proceedings of the Thirteenth*
337 *International Conference on Artificial Intelligence and Statistics*, volume 9 of *Proceedings of*
338 *Machine Learning Research*, pages 249–256, Chia Laguna Resort, Sardinia, Italy, May 2010.
339 PMLR. URL <http://proceedings.mlr.press/v9/glorot10a.html>.
- 340 [9] Elad Hoffer, Itay Hubara, and Daniel Soudry. Train longer, generalize better: closing the
341 generalization gap in large batch training of neural networks. *arXiv:1705.08741 [cs, stat]*, May
342 2017. URL <http://arxiv.org/abs/1705.08741>. arXiv: 1705.08741.
- 343 [10] Samuel Horváth, Aaron Klein, Peter Richtárik, and Cédric Archambeau. Hyperparameter
344 transfer learning with adaptive complexity. *CoRR*, abs/2102.12810, 2021. URL <https://arxiv.org/abs/2102.12810>.
- 346 [11] Xiao Shi Huang and Felipe Pérez. Improving Transformer Optimization Through Better
347 Initialization. page 9.
- 348 [12] Arthur Jacot, Franck Gabriel, and Clément Hongler. Neural Tangent Kernel: Convergence
349 and Generalization in Neural Networks. *arXiv:1806.07572 [cs, math, stat]*, June 2018. URL
350 <http://arxiv.org/abs/1806.07572>. 00000 arXiv: 1806.07572.
- 351 [13] Kevin Jamieson and Ameet Talwalkar. Non-stochastic best arm identification and hyperparame-
352 ter optimization, 2015.
- 353 [14] Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B. Brown, Benjamin Chess, Rewon Child,
354 Scott Gray, Alec Radford, Jeffrey Wu, and Dario Amodei. Scaling Laws for Neural Language
355 Models. *arXiv:2001.08361 [cs, stat]*, January 2020. URL [http://arxiv.org/abs/2001.](http://arxiv.org/abs/2001.08361)
356 [08361](http://arxiv.org/abs/2001.08361). arXiv: 2001.08361.
- 357 [15] Diederik Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint*
358 *arXiv:1412.6980*, 2014.
- 359 [16] Liam Li, Kevin Jamieson, Afshin Rostamizadeh, Ekaterina Gonina, Moritz Hardt, Benjamin
360 Recht, and Ameet Talwalkar. A system for massively parallel hyperparameter tuning, 2020.
- 361 [17] Lisha Li, Kevin Jamieson, Giulia DeSalvo, Afshin Rostamizadeh, and Ameet Talwalkar. Hyper-
362 band: A Novel Bandit-Based Approach to Hyperparameter Optimization. *JMLR 18*, page 52.

- 363 [18] Liyuan Liu, Xiaodong Liu, Jianfeng Gao, Weizhu Chen, and Jiawei Han. Understanding
364 the difficulty of training transformers. In *Proceedings of the 2020 Conference on Empirical*
365 *Methods in Natural Language Processing (EMNLP)*, pages 5747–5763, Online, November
366 2020. Association for Computational Linguistics. doi: 10.18653/v1/2020.emnlp-main.463.
367 URL <https://www.aclweb.org/anthology/2020.emnlp-main.463>.
- 368 [19] Liyuan Liu, Xiaodong Liu, Jianfeng Gao, Weizhu Chen, and Jiawei Han. Understanding the
369 Difficulty of Training Transformers. *arXiv:2004.08249 [cs, stat]*, September 2020. URL
370 <http://arxiv.org/abs/2004.08249>. arXiv: 2004.08249.
- 371 [20] Xiaodong Liu, Pengcheng He, Weizhu Chen, and Jianfeng Gao. Multi-task deep neural
372 networks for natural language understanding. In *Proceedings of the 57th Annual Meeting of*
373 *the Association for Computational Linguistics*, pages 4487–4496, Florence, Italy, July 2019.
374 Association for Computational Linguistics. URL [https://www.aclweb.org/anthology/](https://www.aclweb.org/anthology/P19-1441)
375 [P19-1441](https://www.aclweb.org/anthology/P19-1441).
- 376 [21] Sam McCandlish, Jared Kaplan, Dario Amodei, and OpenAI Dota Team. An Empirical
377 Model of Large-Batch Training. *arXiv:1812.06162 [cs, stat]*, December 2018. URL [http://](http://arxiv.org/abs/1812.06162)
378 arxiv.org/abs/1812.06162. arXiv: 1812.06162.
- 379 [22] Myle Ott, Sergey Edunov, Alexei Baevski, Angela Fan, Sam Gross, Nathan Ng, David Grangier,
380 and Michael Auli. fairseq: A fast, extensible toolkit for sequence modeling, mit license. In
381 *Proceedings of NAACL-HLT 2019: Demonstrations*, 2019.
- 382 [23] Daniel S. Park, Jascha Sohl-Dickstein, Quoc V. Le, and Samuel L. Smith. The Effect of Network
383 Width on Stochastic Gradient Descent and Generalization: an Empirical Study. May 2019.
384 URL <https://arxiv.org/abs/1905.03776v1>.
- 385 [24] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory
386 Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmai-
387 son, Andreas Kopf, Edward Yang, Zachary DeVito, Martin Raison, Alykhan Tejani,
388 Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala.
389 Pytorch: An imperative style, high-performance deep learning library, BSD-style li-
390 cense. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d’Álché Buc, E. Fox, and
391 R. Garnett, editors, *Advances in Neural Information Processing Systems 32*, pages
392 8024–8035. Curran Associates, Inc., 2019. URL [http://papers.neurips.cc/paper/](http://papers.neurips.cc/paper/9015-pytorch-an-imperative-style-high-performance-deep-learning-library.pdf)
393 [9015-pytorch-an-imperative-style-high-performance-deep-learning-library.](http://papers.neurips.cc/paper/9015-pytorch-an-imperative-style-high-performance-deep-learning-library.pdf)
394 [pdf](http://papers.neurips.cc/paper/9015-pytorch-an-imperative-style-high-performance-deep-learning-library.pdf).
- 395 [25] Valerio Perrone, Rodolphe Jenatton, Matthias W Seeger, and Cedric Archambeau. Scalable
396 Hyperparameter Transfer Learning. *NeurIPS 2018*, page 11.
- 397 [26] Martin Popel and Ondřej Bojar. Training Tips for the Transformer Model. *The Prague Bulletin*
398 *of Mathematical Linguistics*, 110(1):43–70, April 2018. ISSN 1804-0462. doi: 10.2478/
399 pralin-2018-0002. URL [http://content.sciendo.com/view/journals/pralin/110/](http://content.sciendo.com/view/journals/pralin/110/1/article-p43.xml)
400 [1/article-p43.xml](http://content.sciendo.com/view/journals/pralin/110/1/article-p43.xml).
- 401 [27] Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena,
402 Yanqi Zhou, Wei Li, and Peter J. Liu. Exploring the Limits of Transfer Learning with a Unified
403 Text-to-Text Transformer. *arXiv:1910.10683 [cs, stat]*, July 2020. URL [http://arxiv.org/](http://arxiv.org/abs/1910.10683)
404 [abs/1910.10683](http://arxiv.org/abs/1910.10683). arXiv: 1910.10683.
- 405 [28] Samuel S. Schoenholz, Justin Gilmer, Surya Ganguli, and Jascha Sohl-Dickstein. Deep Infor-
406 mation Propagation. *arXiv:1611.01232 [cs, stat]*, November 2016. URL [http://arxiv.org/](http://arxiv.org/abs/1611.01232)
407 [abs/1611.01232](http://arxiv.org/abs/1611.01232). arXiv: 1611.01232.
- 408 [29] Christopher J. Shallue, Jaehoon Lee, Joseph Antognini, Jascha Sohl-Dickstein, Roy Frostig,
409 and George E. Dahl. Measuring the Effects of Data Parallelism on Neural Network Training.
410 *arXiv:1811.03600 [cs, stat]*, November 2018. URL <http://arxiv.org/abs/1811.03600>.
411 arXiv: 1811.03600.
- 412 [30] Mohammad Shoeybi, Mostofa Patwary, Raul Puri, Patrick LeGresley, Jared Casper, and Bryan
413 Catanzaro. Megatron-lm: Training multi-billion parameter language models using model
414 parallelism. *CoRR*, abs/1909.08053, 2019. URL <http://arxiv.org/abs/1909.08053>.

- 415 [31] Samuel L. Smith, Pieter-Jan Kindermans, and Quoc V. Le. Don't Decay the Learning Rate,
416 Increase the Batch Size. *arXiv:1711.00489 [cs, stat]*, November 2017. URL <http://arxiv.org/abs/1711.00489>. arXiv: 1711.00489.
- 418 [32] Jasper Snoek, Hugo Larochelle, and Ryan P. Adams. Practical bayesian optimization of machine
419 learning algorithms, 2012.
- 420 [33] Jasper Snoek, Oren Rippel, Kevin Swersky, Ryan Kiros, Nadathur Satish, Narayanan Sundaram,
421 Md. Mostofa Ali Patwary, Prabhath, and Ryan P. Adams. Scalable bayesian optimization using
422 deep neural networks, 2015.
- 423 [34] Danny Stoll, Jörg K.H. Franke, Diane Wagner, Simon Selg, and Frank Hutter. Hyperparameter
424 transfer across developer adjustments, 2021. URL [https://openreview.net/forum?id=](https://openreview.net/forum?id=WP00vDYLXem)
425 [WP00vDYLXem](https://openreview.net/forum?id=WP00vDYLXem).
- 426 [35] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez,
427 Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. *CoRR*, abs/1706.03762, 2017.
428 URL <http://arxiv.org/abs/1706.03762>.
- 429 [36] Alex Wang, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel R Bowman.
430 Glue: A multi-task benchmark and analysis platform for natural language understanding.
431 *EMNLP 2018*, page 353, 2018.
- 432 [37] Alex Wang, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel R. Bowman.
433 Glue: A multi-task benchmark and analysis platform for natural language understanding, 2019.
- 434 [38] Adina Williams, Nikita Nangia, and Samuel Bowman. A broad-coverage challenge corpus
435 for sentence understanding through inference. In *Proceedings of the 2018 Conference of the*
436 *North American Chapter of the Association for Computational Linguistics: Human Language*
437 *Technologies, Volume 1 (Long Papers)*, pages 1112–1122. Association for Computational
438 Linguistics, 2018. URL <http://aclweb.org/anthology/N18-1101>.
- 439 [39] Greg Yang. Tensor Programs III: Neural Matrix Laws. *arXiv:2009.10685 [cs, math]*, September
440 2020. URL <http://arxiv.org/abs/2009.10685>.
- 441 [40] Greg Yang and Edward J. Hu. Feature learning in infinite-width neural networks. *arXiv*, 2020.
- 442 [41] Greg Yang and Sam S. Schoenholz. Deep Mean Field Theory: Layerwise Variance and
443 Width Variation as Methods to Control Gradient Explosion. February 2018. URL <https://openreview.net/forum?id=rJGY8GbR->
444 [.](https://openreview.net/forum?id=rJGY8GbR-)
- 445 [42] Greg Yang and Samuel S. Schoenholz. Mean Field Residual Networks: On the Edge of Chaos.
446 *arXiv:1712.08969 [cond-mat, physics:nlin]*, December 2017. URL [http://arxiv.org/abs/](http://arxiv.org/abs/1712.08969)
447 [1712.08969](http://arxiv.org/abs/1712.08969). 00000 arXiv: 1712.08969.
- 448 [43] Dani Yogatama and Gideon Mann. Efficient Transfer Learning Method for Automatic Hyperpa-
449 rameter Tuning. In *Artificial Intelligence and Statistics*, pages 1077–1085. PMLR, April 2014.
450 URL <http://proceedings.mlr.press/v33/yogatama14.html>. ISSN: 1938-7228.
- 451 [44] Sergey Zagoruyko and Nikos Komodakis. Wide residual networks, 2017.
- 452 [45] Hongyi Zhang, Yann N. Dauphin, and Tengyu Ma. Residual Learning Without Normalization
453 via Better Initialization. In *International Conference on Learning Representations*, 2019. URL
454 <https://openreview.net/forum?id=H1gsz30cKX>.

455 Checklist

456 The checklist follows the references. Please read the checklist guidelines carefully for information on
457 how to answer these questions. For each question, change the default **[TODO]** to **[Yes]**, **[No]**, or
458 **[N/A]**. You are strongly encouraged to include a **justification to your answer**, either by referencing
459 the appropriate section of your paper or providing a brief inline description. For example:

- 460 • Did you include the license to the code and datasets? **[Yes]** See Section ?.
- 461 • Did you include the license to the code and datasets? **[No]** The code and the data are
462 proprietary.
- 463 • Did you include the license to the code and datasets? **[N/A]**

464 Please do not modify the questions and only use the provided macros for your answers. Note that the
465 Checklist section does not count towards the page limit. In your paper, please delete this instructions
466 block and only keep the Checklist section heading above along with the questions/answers below.

- 467 1. For all authors...
 - 468 (a) Do the main claims made in the abstract and introduction accurately reflect the paper's
469 contributions and scope? **[Yes]**. See abstract and Section 1.
 - 470 (b) Did you describe the limitations of your work? **[Yes]**. See Section 5.1 and Section 8.
 - 471 (c) Did you discuss any potential negative societal impacts of your work? **[Yes]**. See
472 Section 8.
 - 473 (d) Have you read the ethics review guidelines and ensured that your paper conforms to
474 them? **[Yes]**.
- 475 2. If you are including theoretical results...
 - 476 (a) Did you state the full set of assumptions of all theoretical results? **[N/A]**
 - 477 (b) Did you include complete proofs of all theoretical results? **[N/A]**
- 478 3. If you ran experiments...
 - 479 (a) Did you include the code, data, and instructions needed to reproduce the main exper-
480 imental results (either in the supplemental material or as a URL)? **[No]** Our main
481 experiments are quite expensive and depend on details of our internal cluster, so we do
482 not provide code to exactly reproduce it. However we include a package so that anyone
483 can use μ Transfer and describe all the details to reproduce it.
 - 484 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
485 were chosen)? **[Yes]**. See Section 6 and Appendix C.
 - 486 (c) Did you report error bars (e.g., with respect to the random seed after running experi-
487 ments multiple times)? **[Yes]**
 - 488 (d) Did you include the total amount of compute and the type of resources used (e.g., type
489 of GPUs, internal cluster, or cloud provider)? **[Yes]**. See Section 6.
- 490 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
 - 491 (a) If your work uses existing assets, did you cite the creators? **[Yes]**. We used Megatron,
492 fairseq, pytorch and cited all of them.
 - 493 (b) Did you mention the license of the assets? **[Yes]**. In the references.
 - 494 (c) Did you include any new assets either in the supplemental material or as a URL? **[Yes]**
495 . Yes, we are releasing a new Pytorch package.
 - 496 (d) Did you discuss whether and how consent was obtained from people whose data you're
497 using/curating? **[N/A]**
 - 498 (e) Did you discuss whether the data you are using/curating contains personally identifiable
499 information or offensive content? **[N/A]**
- 500 5. If you used crowdsourcing or conducted research with human subjects...
 - 501 (a) Did you include the full text of instructions given to participants and screenshots, if
502 applicable? **[N/A]**
 - 503 (b) Did you describe any potential participant risks, with links to Institutional Review
504 Board (IRB) approvals, if applicable? **[N/A]**
 - 505 (c) Did you include the estimated hourly wage paid to participants and the total amount
506 spent on participant compensation? **[N/A]**

507 A Detailed Discussions on Related Works

508 A.1 Hyperparameter Tuning

509 Many have sought to speedup hyperparameter tuning beyond the simple grid or random search. Snoek
510 et al. [32] treated hyperparameter tuning as an optimization process and used Bayesian optimization
511 by treating the performance of each hyperparameter combination as a sample from a Gaussian process
512 (GP). Snoek et al. [33] further improved the runtime by swapping the GP with a neural network.
513 Another thread of work investigated how massively parallel infrasture can be used for efficient tuning
514 under the multi-arm bandit problem [13, 16]. There are also dedicated tools such as Optuna [4] and
515 Talos [3] which integrate with existing deep learning frameworks and provide an easy way to apply
516 more advanced tuning techniques.

517 Our approach is distinct from all of the above in that it does not work on the hyperparameter
518 optimization process itself. Instead, it decouples the size of the target model from the tuning
519 cost, which was not feasible prior to this work. This means that **no matter how large the target
520 model is, we can always use a fixed-sized proxy model to probe its hyperparameter landscape**
521 Nevertheless, our method is complementary, as the above approaches can naturally be applied to the
522 tuning of the proxy model; it is only for scientific reasons that we use either grid search or random
523 search throughout this work.

524 A.2 Previously Proposed Scaling Rules of Hyperparameters

525 **(Learning Rate, Batch Size) Scaling** [31] proposed to scale learning rate with batch size while
526 fixing the total epochs of training; [9] proposed to scale learning rate as $\sqrt{batchsize}$ while fixing
527 the total number of steps of training. However, [29] showed that there’s no consistent (learning
528 rate, batch size) scaling law across a range of dataset and models. Later, [21] studied the trade-off
529 of training steps vs computation as a result of changing batch size. They proposed an equation of
530 $a/(1 + b/batchsize)$, where a and b are task- and model-specific constants, for the optimal learning
531 rate (see their fig 3 and fig 5). This law suggests that for sufficiently large batch size, the optimal
532 learning rate is roughly constant.¹⁵ This supports our results here as well as the empirical results in
533 [29, fig 8].

534 **Learning Rate Scaling with Width** Assuming that the optimal learning rate should scale with
535 batch size following [31], [23] empirically investigated how the “noise ratio” $LR/batchsize$ scales
536 with width for MLP and CNNs in NTK parametrization (NTP) or standard parametrization (NTP)
537 trained with SGD. They claimed that, in networks without batch normalization, the optimal noise
538 ratio is constant in SP but scales like $1/width$ for NTP. However, they found this law breaks down
539 for networks with normalization.

540 Here in our work, Fig. 3 contradicts their results on SP MLP by showing the optimal learning rate
541 (fixing batch size) shifts with width. We believe this difference is 1) due to their erroneous assumption
542 that optimal learning rate scales with batch size (as debunked by [21, 29]) and 2) because their SP
543 experiments were done by fixing the learning rate and only sweeping batch size.

544 Furthermore, Fig. 1 clearly shows the optimal learning rate is *not* constant in SP for transformers
545 (trained with Adam). Other differences in our works include our applicability to 1) networks with
546 normalization, 2) Adam and other adaptive optimizers, 3) our empirical validation of transfer across
547 depth and sequence length, and 4) explicit validation of tuning via hyperparameter transfer on large
548 models like BERT-large.

549 Finally, as argued in [40] and Appendix G.3, SP and NTP lead to bad infinite-width limits in contrast
550 to μP and hence are suboptimal for wide neural networks. For example, sufficiently wide neural
551 networks in SP and NTP would lose the ability to learn features, as concretely demonstrated on
552 word2vec in [40].

553 **Input Layer Parametrization** While typically, the input layer is initialized with fanin initialization,
554 in language models where the input and output layers are shared (corresponding to word embeddings),
555 it can actually be more natural to use a fanout initialization (corresponding to fanin initialization of

¹⁵while the optimal learning is roughly linear in batch size when the latter is small

556 the output layer). In fact, we found that fairseq [22] by default actually implements our proposed
557 input layer parametrization (both the fanout initialization and the $\sqrt{\text{fan_out}}$ multiplier).¹⁶

558 **From the Theory of Infinite-Width to the Practice of Finite-Width Neural Networks and Back**
559 [40] introduced μP as the unique parametrization that enables all layers of a neural network to learn
560 features in the infinite-width limit, especially in contrast to the NTK parametrization [12] (which
561 gives rise to the NTK limit) that does not learn features in the limit. Based on this theoretical
562 insight, in Appendix G.3, we argue that μP should also be the unique parametrization that allows
563 hyperparameter transfer across width; in short this is because it both 1) preserves feature learning, so
564 that performance on feature learning tasks (such as language model pretraining) does not become
565 trivial in the limit, and 2) ensures each parameter tensor is not stuck at initialization in the large
566 width limit, so that its learning rate does not become meaningless. At the same time, our results
567 here suggest that μP is indeed the *correct* parametrization for large neural networks and thus provide
568 empirical motivation for the theoretical study of the infinite-width μP limit.

569 B Which Hyperparameters Can Be Transferred? (Continued)

570 B.1 Further Discussions on Hyperparameter Categories

571 Below, we discuss the reasoning behind each kind, which are supported by our empirical evidence
572 collected in Fig. 4 on Transformers as well as those in Appendix D.1 on ResNet.

573 **Transferable Hyperparameters** In Table 2, we summarize which hyperparameters can be trans-
574 ferred across training scale. The transfer across *width*, as explained in Section 2, is theoretically
575 justified, while we present the transfer across the other dimensions as empirical results.

576 These cover most of the well-known and important hyperparameters when the need for regularization
577 is not paramount, e.g., during large scale language model pretraining. Parameter Multipliers are not
578 well-known hyperparameters, yet we include them here as they serve a bridge between SP and μP and
579 can impact model performance in practice. Concretely, any SP and μP neural networks of the same
580 width can have their Parameter Multipliers tuned so that their training dynamics become identical.

581 **Hyperparameters That Don’t Transfer Well** Not all hyperparameters transfer well even if we
582 use μP . In particular, those whose primary function is to regularize training to mitigate “overfitting”
583 tend not to transfer well. Indeed, intuitively, regularization needs to be applied more heavily in larger
584 models, so naturally we do not expect the same regularization hyperparameters to stay constant across
585 model sizes.

586 To the best of our knowledge, there is no strict separation between hyperparameters that regularize
587 and those that don’t. However, conventional wisdom tells us that there exists a spectrum of how
588 much regularizing effect a hyperparameter has. For example, dropout probability and weight decay
589 are among those whose primary function is to regularize, whereas batch size and learning rate might
590 regularize training in some cases but affect the dynamics more so in other ways. Our empirical
591 exploration tells us that the former do not transfer well, while the latter do. Our subsequent discussion
592 will focus on the latter; we leave to future works the expansion to the former.

593 **Hyperparameters Transferred Across** We have left out a category of hyperparameters that defines
594 the training *scale*, or in practical terms, training cost. This includes 1) those that define how many
595 operations a model’s forward/backward pass takes, such as the model’s width, depth, and in the case
596 of language modeling, sequence length; and 2) those that define how many such passes are performed,
597 such as batch size and number of training steps.

598 As recent works have shown [6, 14], improvements along any of these *scale* dimensions lead to
599 apparently sustainable gain in performance; as a result, we are primarily interested in transferring
600 other hyperparameters *across* these dimensions that define scale, rather than finding the optimal
601 scale.¹⁷ This category of hyperparameters is particularly crucial as one can speedup training by

¹⁶But it certainly does not implement other parts of our parametrization, like Adam learning rate scaling or the output multiplier.

¹⁷In particular, we are not fixing the total training FLOPs when we scale, which requires understanding the tradeoff of different scale hyperparameters. For example, when we transfer across batch size, we *fix* the number of steps of training (*not* the number of epochs), so that the total FLOPs scales linearly.

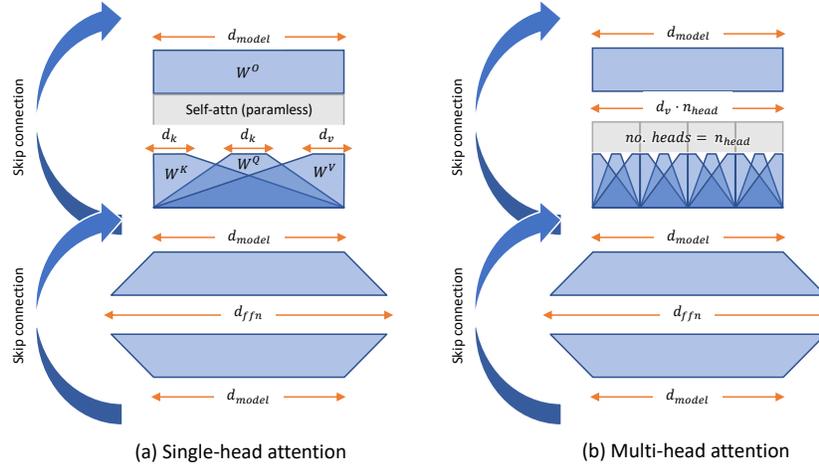


Figure 6: Schematics of each Transformer layer. Commonly, the key and value dimensions d_k and d_v are both set to d_{model}/n_{head} , and this is referred to as d_{head} .

602 downsizing in one or multiple such dimensions. Indeed, it’s very common for practitioners to
 603 implicitly transfer hyperparameters across the number of training samples by tuning on only a subset
 604 of the full training data.

605 Our insights from the infinite-width limit inspired us to explore hyperparameter transfer across *width*,
 606 which does not work under SP as we have shown earlier. Building upon our success with width,
 607 which is well explained theoretically, we hope to push the limit of compute-saving by investigating
 608 the other dimensions empirically. To the best of our knowledge, the transferability of optimal
 609 hyperparameters across depth, batch size, sequence length, and training time has not been rigorously
 610 investigated previously, with the main exception of the literature on (learning rate, batch size) scaling
 611 [29, 31] where our transferability result of learning rate across batch size recapitulates [21].¹⁸ See
 612 Appendix A.2 on how our results relate to prior works. We will primarily focus on the Transformer
 613 architecture in the main text with evidence for ResNet in Appendix D.1.

614 B.2 On the Definitions of Width

615 Our theory allows more general notions of width. This is especially relevant in Transformers, where
 616 $d_{model}, d_{head} = d_k, d_v, n_{head}, d_{ffn}$ (see Fig. 6) can all be construed as measures of width. We briefly
 617 discuss these here, with more theoretical justification in Appendix G.2.1 and empirical validation
 618 below.

619 **Varying Width Ratio** So far we have assumed that every hidden layer is widened by the same
 620 factor. But in fact we can widen different hidden layers differently. This is useful, for example, in a
 621 Transformer where we may want to use a smaller d_{ffn} during tuning. If we are using Adam, as long
 622 as the width of every layer still tends to infinity, we still obtain approximately the same limit¹⁹, so the
 623 hyperparameter transfer remains theoretically justified.

624 See Fig. 7 for an empirical validation on IWSLT-14 using a Transformer.

625 **Number of Attention Heads** In attention-based models, one typically splits hidden size into
 626 multiple attention heads following $d_{model} = d_{head} \times n_{heads}$. So far we have assumed d_{head} and
 627 d_{model} to be width, but it’s possible and potentially advantageous to fix d_{head} and treat n_{heads} as
 628 the width, or increasing both simultaneously. This allows our technique to handle many popular
 629 models, including GPT-3 [6], which scale up by fixing d_{head} and increasing n_{head} . See Fig. 8 for an
 630 empirical validation on Wikitext-2.

¹⁸There’s also a literature on the proper initialization for training deep networks effectively (e.g. [5, 11, 19, 28, 41, 42, 45]), but they do not study the *transferability* per se. See Appendix A.2

¹⁹This also applies for SGD, but we need more involved scaling to keep the limit approximately the same.

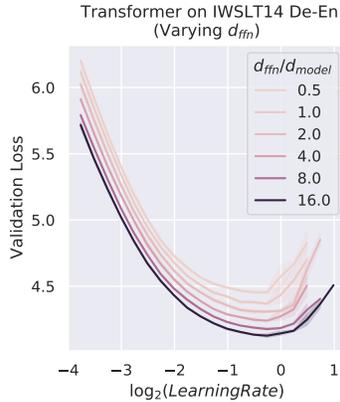


Figure 7: Learning rate landscape in μP is stable even if we vary d_{ffn} by a factor of 32, fixing d_{model} .

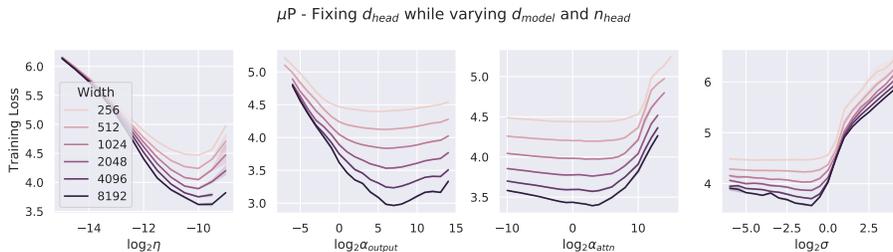


Figure 8: Hyperparameters transfer across width when we fix d_{head} and vary d_{model} and n_{head} . α_{output} , α_{attn} are multipliers for output and key weights, and σ is initialization standard deviation.

631 **Varying Just the Width of Attention Heads** A specific useful instance of varying width ratio is
 632 decoupling the key and value dimensions d_k and d_v and scaling d_k differently from (typically larger
 633 than) d_{model}/n_{heads} . This works as long as we use $1/d$ scaled-attention as in Definition 4.1 (instead
 634 of $1/\sqrt{d}$ as is done commonly). When tuning on the small proxy model, if d_k is too small, the
 635 hyperparameter landscape can be quite noisy. Keeping d_k relatively large while shrinking all other
 636 dimensions solves this problem, while still obtaining significant speedup.

637 C Experimental Details

638 C.1 IWSLT

639 IWSLT14 De-En is a well-known machine translation benchmark. We use a Transformer implemented
 640 in fairseq [22] with a default $d_{model} = 1/4d_{ffn} = 512$ and $d_k = d_q = d_v = d_{model}/n_{heads} = 128$
 641 (amounting to 40M parameters), which we denote as the *1x model*. For transfer, we tune on a proxy
 642 model with the same n_{head} but with d_{model} and other dimensions 4 times smaller; we will call this
 643 the *0.25x model* (but it has 4M parameters). All models are trained with Adam for 100 epochs and
 644 validated at the end of every epoch. We tune via random search the learning rate η , the output layer
 645 parameter multiplier α_{output} , and the attention key-projection weight multiplier α_{attn} following the
 646 grid

- 647 • $\eta: 5 \times 10^{-4} \times 2^z$, where $z \in \{-1.5, -1.25, -1, \dots, 1.25\}$
- 648 • $\alpha_{output}: 2^z$, where $z \in \{-8, -7, -6, \dots, 7\}$
- 649 • $\alpha_{attn}: 2^z$, where $z \in \{-3, -2, -1, \dots, 8\}$

650 C.2 WMT

651 We scale up to WMT14 En-De using the large Transformer from [35], with a $d_{model} = 1/4d_{ffn} =$
 652 1024 and $d_q = d_k = d_v = d_{model}/n_{heads} = 64$. We use the exact same setup and reproduce their
 653 result as our baseline. Then, we build the proxy model by shrinking the target model’s d_{model} from

654 the original 1024 to 256, d_{ffn} from 4096 to 256 and n_{heads} from 16 to 4. This reduces the total
 655 parameter count from 211M to 15M. We then perform the hyperparameter search on the proxy
 656 model and take the best according to validation loss, before testing on the target model. We tune via
 657 random search the learning rate η , the output layer parameter multiplier α_{output} , and the attention
 658 key-projection weight multiplier α_{attn} following the grid

- 659 • $\eta: 6 \times 10^{-4} \times 2^z$, where $z \in \{-1.5, -1.25, -1, \dots, 1.25\}$
- 660 • $\alpha_{output}: 2^z$, where $z \in \{-8, -7, -6, \dots, 7\}$
- 661 • $\alpha_{attn}: 2^z$, where $z \in \{-3, -2, -1, \dots, 8\}$

662 C.3 BERT

663 **Details of BERT Prototype** Our proxy model has 10 Transformer layers with $d_{model} = d_{ffn} =$
 664 256. We also reduce the number of attention heads to 8 with a d_{head} of 32. We call it BERT Prototype
 665 since we can increase its width and depth according to our definitions to recover both BERT Base
 666 and BERT Large, which enables us to sweep hyperparameters once and use for both models. Overall,
 667 BERT Prototype has 13M trainable parameters, a fraction of the 110M in BERT Base and the 350M
 668 in BERT Large.

669 **Hyperparameters Tuned for Pretraining** We tune the following hyperparameters for pretraining:
 670 Adam learning rate η , embedding learning rate η_{emb} , output weight multiplier α_{output} , attention
 671 logits multiplier α_{attn} , layernorm gain multiplier $\alpha_{LN_{gain}}$, and bias multiplier α_{bias} .

672 We sample 256 combinations from the follow grid:

- 673 • $\eta: 1 \times 10^{-4} \times 2^z$, where $z \in \{1.5, 2, 2.5, 3, 3.5\}$
- 674 • $\eta_{emb}: 1 \times 10^{-4} \times 2^z$, where $z \in \{-1, -0.5, 0, 0.5, 1\}$
- 675 • $\alpha_{output}: 2^z$, where $z \in \{2, 4, 6\}$
- 676 • $\alpha_{attn}: 2^z$, where $z \in \{3, 3.5, 4, \dots, 7\}$
- 677 • $\alpha_{LN_{gain}}: 2^z$, where $z \in \{8.5, 9, 9.5, 10, 10.5\}$
- 678 • $\alpha_{bias}: 2^z$, where $z \in \{8.5, 9, 9.5, 10, 10.5\}$

679 The ranges are chosen to include the implicit choices of these hyperparameters in SP BERT Large.

680 **Finetuning Procedure and Hyperparameters** We hand-pick the finetuning hyperparameters after
 681 training the full-sized model. As regularization is an essential ingredient in successful finetuning, we
 682 do not perform hyperparameter transfer (at least the suite of techniques presented in this work) (see
 683 Table 1). We focus on MNLI [38] and QQP, which are two representative tasks from GLUE [36].
 684 Following [20], we used Adam [15] with a learning rate of 5×10^{-5} and a batch size of 64. The
 685 maximum number of epochs was set to 5. A linear learning rate decay schedule with warm-up of 0.1
 686 was used. All the texts were tokenized using wordpieces and were chopped to spans no longer than
 687 128 tokens.

688 D Additional Experiments

689 D.1 Experiments on ResNets

690 D.1.1 ResNet on CIFAR-10

691 **Setup** For this case we use Davidnet [2], a ResNet variant that trains quickly on CIFAR-10, so
 692 as to efficiently investigate its hyperparameter landscape. We train with SGD on CIFAR-10 for 10
 693 epochs; all results are averaged over 15 random seeds. We use a width multiplier to identify models of
 694 different width, and a multiplier of 1 corresponds to the original model in [2]. We look at validation
 695 accuracy here as the model barely overfits, and our observations will hold for the training accuracy as
 696 well. We first conduct a learning rate sweep for models of different widths using SP; the result is
 697 shown in Fig. 9, on the left.

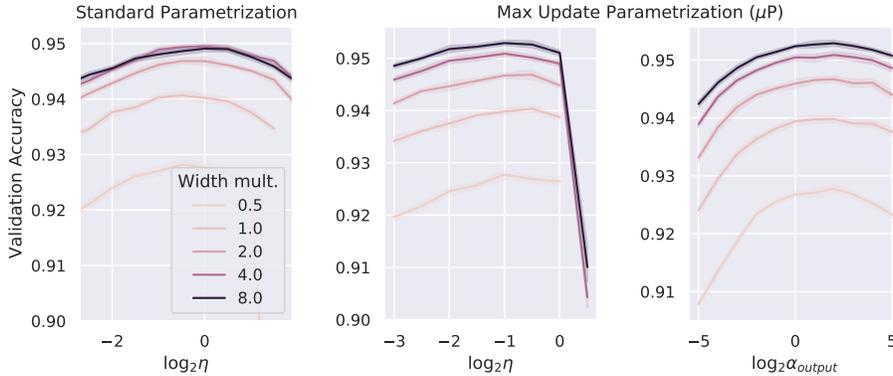


Figure 9: ResNet on CIFAR-10 for different widths (compared to a base network). On the **left**, the widest network SP underperforms; on the **right**, the μP network has a more consistent hyperparameter landscape and performs better. Both networks are tuned at the smallest width for the hyperparameter (η or α_{output}) not in the x-axis.

698 **Hyperparameter Stability** Note that the best model with a width multiplier of 8 underperforms
 699 that with a multiplier of 4. We run the same sweep with μP , along with a sweep of the output
 700 multiplier (α_{output}); the result is shown in Fig. 9, on the right. We notice that wider models always
 701 perform better under μP and that the optimal learning rate η and α_{output} are stable across width.

702 **Hyperparameter Transfer** Next, we perform a grid search for learning rate (η) and α_{output} on
 703 the 0.5x model for both SP and μP .²⁰ Then, we take the best combination and test on the 8x model,
 704 simulating how a practitioner might use μ Transfer. The result is shown in Table 7, where μP
 705 outperforms SP by $0.43\% \pm .001\%$.

Table 7: Transferring the best learning rate (η) and α_{output} from widening factor 0.5 to 8; μP significantly outperforms SP given the same search grid. The best hyperparameters are different as the models are parametrized to be identical at 1x width.²⁰

| Transfer Setup | Best η | Best α_{output} | Valid. Acc. (0.5x) | Valid. Acc. (8x) |
|----------------|-------------|------------------------|--------------------|------------------|
| SP | 0.707 | 4 | 92.82% | 94.86% |
| μP | 0.5 | 4 | 92.78% | 95.29% |

706 D.1.2 Wide ResNet on ImageNet

707 **Setup** For this case we use Wide-Resnet, or WRN [44], a ResNet variant with more channels per
 708 layer, to further showcase hyperparameter transfer across width, i.e., number of channels. We train
 709 with SGD on ImageNet for 50 epochs following standard data augmentation procedures. We use
 710 a width multiplier to identify models of different width, and a multiplier of 1 corresponds to the
 711 original WRN-50-2-bottleneck in [44].

712 **Hyperparameter Transfer** We start with a proxy model with a width multiplier of 0.125 and tune
 713 several hyperparameters using the following grid:

- 714 • $\eta: 1 \times 2.048 \times 2^z$, where $z \in \{-5, -4, -3, \dots, 4\}$
- 715 • $\alpha_{output}: 10 \times 2^z$, where $z \in \{-5, -4, -3, \dots, 4\}$
- 716 • weight decay co-efficient $\gamma: 3.05 \times 10^{-5} \times 2^z$, where $z \in \{-2, -1.5, -1, \dots, 1.5\}$
- 717 • SGD momentum $\beta: 0.875 \times 2^z$, where $z \in \{-2, -1.5, -1, \dots, 1.5\}$

718 The grid is centered around the default hyperparameters used by [1] for ResNet-50; while not expected
 719 to be competitive for WRN, they represent a reasonable starting point for our experiment.

²⁰Here we tune the 0.5x model instead of the 1x model to simulate the situation that one does “exploratory work” on the 1x model but, when scaling up, would like to tune faster by using a smaller proxy model.

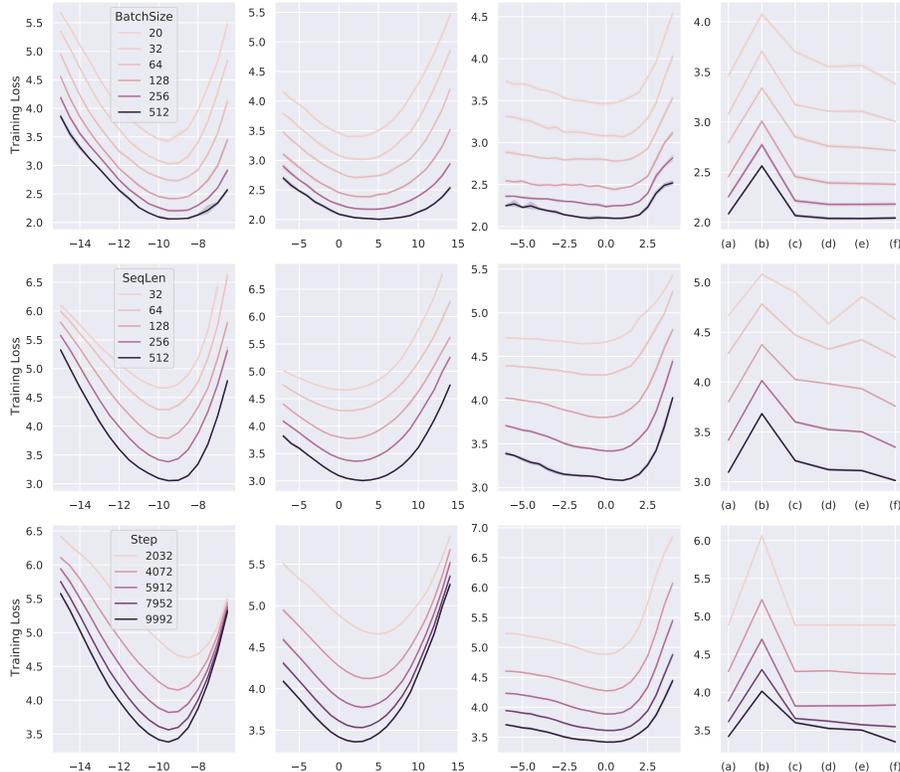


Figure 10: **Empirical validation of Hyperparameter Transfer across Batch Size, Sequence Length, and Training Time on pre-LN Transformers.** Same setting as Fig. 4. Despite some shift, the optimal hyperparameters are roughly stable when transferring from batch size 32, sequence length 128, and 5000 training steps.

720 We randomly sample 64 hyperparameter combinations from the grid and train for 50 epochs, before
 721 selecting the one with the highest top-1 validation accuracy. Then, we scale up the model following
 722 both μ P and SP and run with the same hyperparameters we just selected. The result is shown in
 723 Table 8, where μ P outperforms SP by 0.41% in terms of top-1 validation accuracy.

Table 8: Transferring the best learning rate (η), α_{output} , γ , and β from widening factor 0.125 to 1; μ P significantly outperforms SP given the same search grid.

| Transfer Setup | Best η | Best α_{output} | Best γ | Best β | Valid. Acc. (0.125x) | Valid. Acc. (1x) |
|----------------|-------------|------------------------|---------------|--------------|----------------------|------------------|
| SP | 32.768 | .625 | .000015 | .4375 | 58.12% | 76.75% |
| μ P | 32.768 | .625 | .000015 | .4375 | 58.12% | 77.16% |

724 D.2 Experiments on Transformers

725 D.2.1 Verifying Transfer across Batch Size, Sequence Length, and Training Time on 726 Wikitext-2

727 See Fig. 10.

728 D.3 Post-Layernorm Transformers

729 Fig. 11 shows the transferability of learning rate, α_{output} , initialization standard deviation, and Adam
 730 β_2 across width, batch size, sequence length, and training steps for post-layernorm transformers.
 731 However, in general, we find transfer across depth to be fragile.

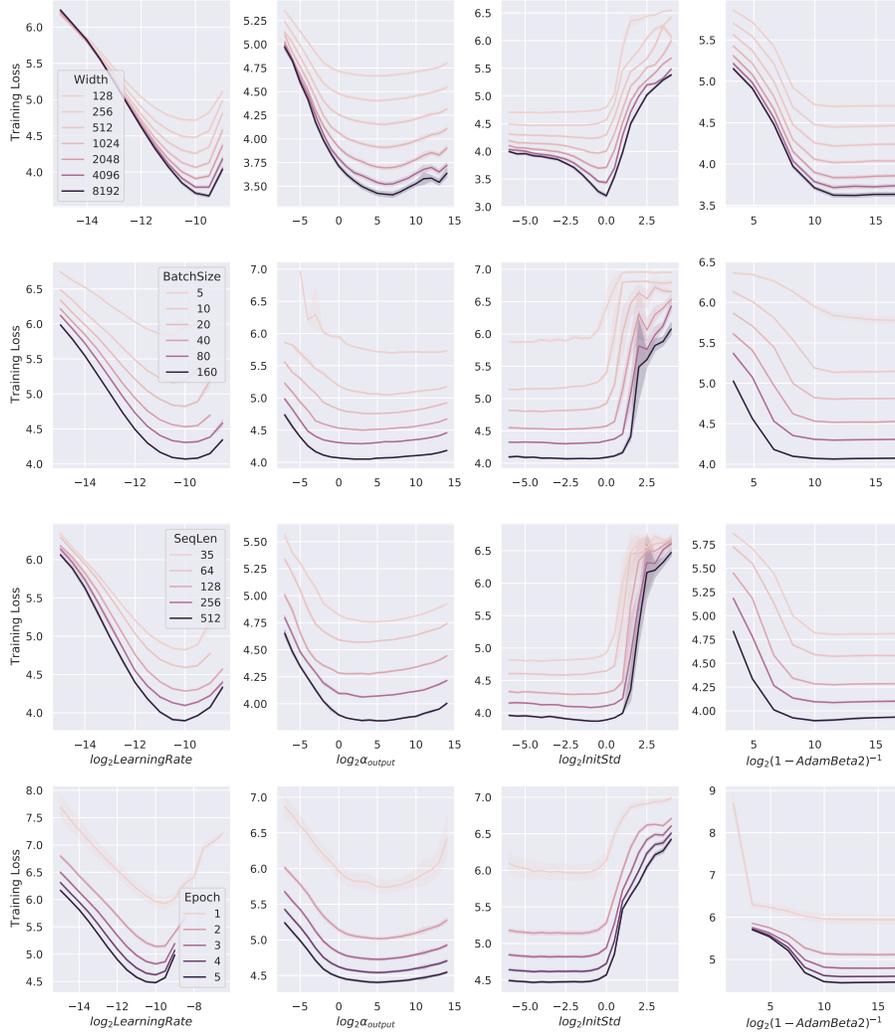


Figure 11: **Empirical validation of Hyperparameter Transfer for Post-LN Transformers.** Same setting as Fig. 4.

732 D.3.1 Hyperparameter Instability of SP Transformers

733 Fig. 12 and Fig. 13 show the hyperparameter instability inherent in SP Transformers.

734 E Implementing Hyperparameter Transfer in a Jiffy

735 As we have shown, one can enable hyperparameter transfer by just reparametrizing the desired
 736 model in Maximal Update Parametrization (μ P). While conceptually simple, switching from Standard
 737 Parametrization (SP) to μ P can be error-prone, as popular deep learning frameworks are built around
 738 SP. We strive to build a tool that fulfills two goals:

- 739 1. Minimize code changes when switching to μ P;
- 740 2. Keep model behavior invariant, under this switch, at a given model `base_width`.

741 The latter goal, which we call *parametrization backward compatibility*, ensures that any code base
 742 works exactly as before when model width equals `base_width`, similar to Eq. (4), e.g. the loss at
 743 any time step remains exactly the same before and after the switch to μ P. Of course, when width
 744 differs from `base_width`, the model behavior necessarily changes so that hyperparameters can be
 745 transferred. Most commonly, the user should set the `base_width` to be the width of the target model

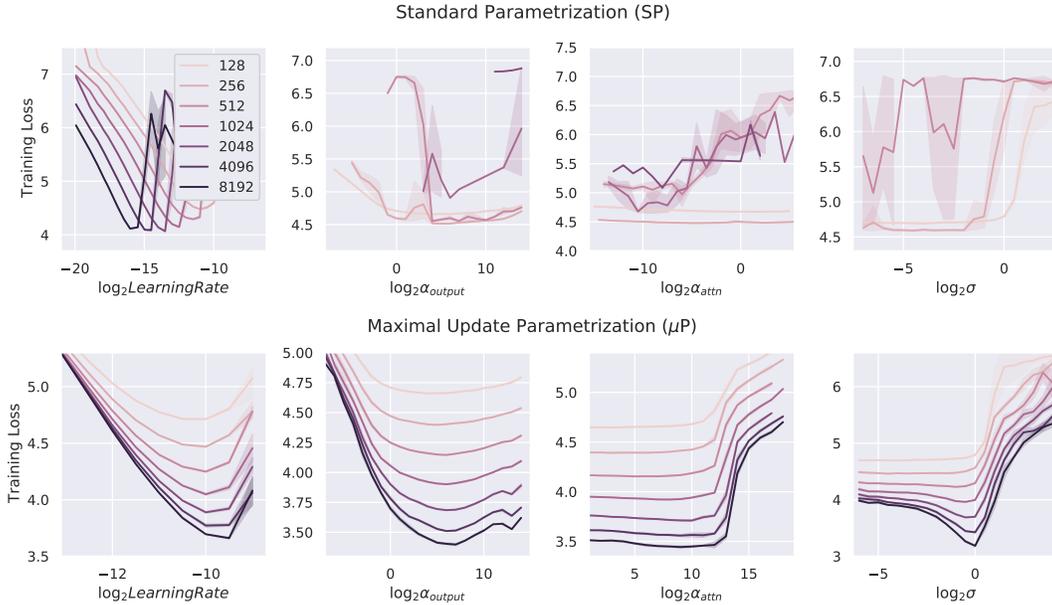


Figure 12: Post-layernorm Transformer with SP and μP on Wikitext-2. We sweep one hyperparameter across width (d_{model}) at a time while keeping the rest fixed; we also scale d_{head} linearly with d_{model} and fixing n_{heads} . $\alpha_{output}, \alpha_{attn}$ are multipliers for output and key weights, and σ is initialization standard deviation. This yields unstable result for SP, as expected, where missing points/curves represent divergence; in μP , the optimal hyperparameter choices stabilize as width increases.

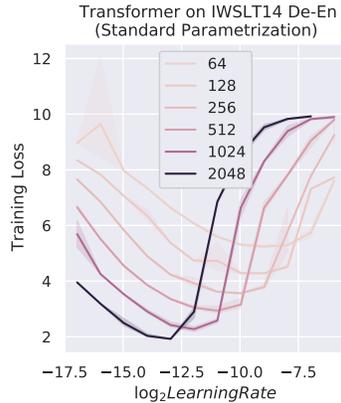


Figure 13: Learning rate landscape is highly unstable under standard parametrization in IWSLT.

746 (e.g. BERT-large or T5-large). Then one can tune a proxy model with e.g. $width = base_width/4$
 747 to obtain the optimal hyperparameters for the target model. In addition, if one wishes to scale up
 748 further e.g. $width = 4 \times base_width$, then these hyperparameters remain optimal. Of course, depth,
 749 batch size, and sequence lengths can be scaled up and down as well according to Fig. 10.

750 **The MUP Package** We provide our tool as a Python package called MUP designed to work with
 751 PyTorch. For the most generic use case, where one scales the widths of all layers at once, the
 752 transition to μP boils down to 3 steps:

- 753 1. Replace the input and the output layers with counterparts in `MUP.Layer` and specify a
- 754 `base_width` for both;
- 755 2. Ensure all other layers are initialized with `fan_in` initialization;²¹

²¹This is the default behavior for Pytorch `nn.Linear` layers, but some code bases then manually overrides this initialization e.g. with constant init.

756

3. Replace the optimizer (e.g., Adam) with the counterpart (e.g., MuAdam) in `MUP.optim`.

757

What Happens in the MUP Package The MuLayers take care of the parametrization for input/output layers and label them for the optimizer. The MuOptimizer reads the `base_width` from the input layer and calculates the learning rate scaling for all parameters. For example, MuAdam scales the learning rate for the input/output layers like $\Theta(1/\sqrt{width})$, one-dimensional parameters (e.g., gains and biases) like $\Theta(1)$, and other parameters like $\Theta(1/width)$. It might seem odd that the optimizer “micromanages” the learning rate for gains and biases; this design choice is motivated by minimizing the required code change by the user, as the alternative is to replace every layer that has gains or biases.

765

F Width-Related Training Issues are Hyperparameter Issues

766

Large transformers are famously fickle to train [18, 26]. This is especially true in low-precision floating point formats such as `float16` that is required to train enormous models today. Throughout our investigations, we find that, with the hyperparameters transferred from the proxy model (which should be close to optimal), we do not observe such instability. This suggests that the training instability problems could in fact be hyperparameter problems, not necessarily structural issues associated with the architecture or optimizer (which work just fine when tuned properly).²² We hypothesize that

773

The common manifestations of training instability are caused by practitioners naively transferring learning rate and other hyperparameters tuned on a small transformer.

774

This is certainly consistent with Fig. 1, which shows that the optimal learning rate for small transformers can lead to trivial performance in large transformers. We support this hypothesis further by *reverse-transferring the instability-inducing hyperparameters from a large transformer to a small one and replicating the training instability*. This is shown in Fig. 14.

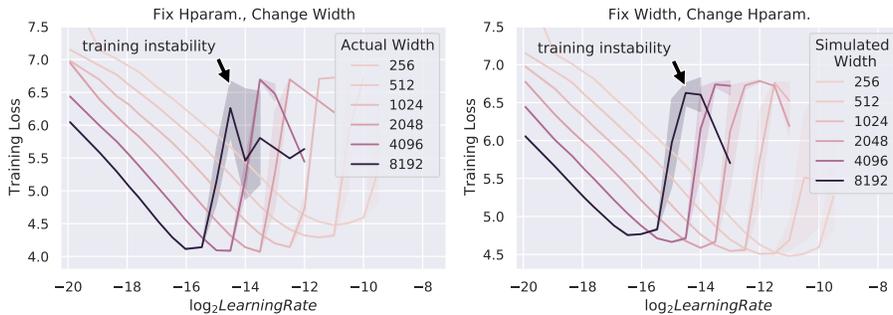


Figure 14: **Replicating training instability on a small transformer by reverse-transferring hyperparameters.** These experiments concern 2-layer Transformers in Standard Parametrization (SP) on Wikitext-2, trained with Adam, where width is defined as $d_{model} = d_{ffn}$. (Left) LR-vs-loss for wider and wider transformers. (Right) Likewise for *simulated width*: Here each point $(\log_2 \eta, loss)$ for simulated width n indicates the loss from training a width-256 μP Transformer with base width n and LR η (i.e. loosely speaking, it’s using LR transferred from η in a width- n SP Transformer). **Takeaway:** The overall shapes of the curves are identical between the left and right plots²³; in particular, a learning rate leads to instability in a wide model iff it does so when transferred back to a narrow model.

²²Of course, it is still worthwhile to research layers or architectures that are more insensitive to poor hyperparameter choices. But with our results, good hyperparameters are obtained much more easily, so this issue is now much less urgent.

²³Note that the curves on the left are “lower” than curves on the right. This just reflects the increasing capacity of wider models able to fit the training data better, so is orthogonal to our point.

Table 9: Expected entry size of Av for different matrices A and vector v correlated with each other, both having entries of size $\Theta(1)$.

| | Standard Gaussian $A \in \mathbb{R}^{n \times n}$ | (Nonlinear) Tensor Product $A \in \mathbb{R}^{n \times n}$ | Vector $A \in \mathbb{R}^{1 \times n}$ |
|--------------------|--|---|---|
| Entry size of Av | $\Theta(\sqrt{n})$ | $\Theta(n)$ | $\Theta(n)$ |

778 G An Intuitive Introduction to the Theory of Maximal Update 779 Parametrization

780 In what follows, we seek to describe useful intuitions and rule of thumbs that would be helpful
781 to practitioners and empirical researchers alike in figuring out what is the right neural network
782 parametrization. Readers needing rigorous justification can generally find it in [39, 40].

783 G.1 Behaviors of Gaussian Matrices vs Tensor Product Matrices

784 Central to the derivation of μP for any architecture are key insights on the behaviors of two kinds of
785 random matrices: 1) iid Gaussian random matrix and 2) tensor product matrix (by which we mean a
786 sum of outer products) and more generally what we call *nonlinear* tensor product matrix (see Eq. (6)).
787 For example, a neural network, randomly initialized in the typical way, will have each weight matrix
788 look like the former. However, every step of training by gradient descent adds a sum of outer products
789 to this initial matrix, so that the *change in weights* constitute a tensor product matrix. For Adam,
790 the change in weights is not a tensor product but a more general *nonlinear tensor product matrix*
791 (see Eq. (6)). In this section, we will particularly focus on the *right scaling* for the entries of such
792 matrices, leading to a discussion of the *right neural network parametrization* in the next section. We
793 concentrate on the key heuristics but eschew burdensome rigor.

794 **Key Insights** Consider a random vector $v \in \mathbb{R}^n$ with approximately iid entries and a random
795 matrix A of either size $n \times n$ or $1 \times n$, both having entries of size $\Theta(1)$.²⁴ In the context of deep
796 learning, v for example can be an activation vector in an MLP, a Gaussian A the hidden weights at
797 initialization, a (nonlinear) tensor product A the change in hidden weights due to training, and a
798 vector A the readout layer weights. Then Av corresponds to a part of the next layer preactivation
799 or the network output. To make sure the preactivations and the output don't blow up, we thus need
800 to understand the scale of Av , especially in the general case where A is correlated with v .²⁵ This is
801 summarized in Table 9, with the derivations below. Intuitively, a (nonlinear) tensor product or vector
802 A will interact with a correlated v via Law of Large Numbers, hence the n -scaling, while a Gaussian
803 A interacts with v via Central Limit Theorem, hence the \sqrt{n} -scaling.

804 In the derivations below, we answer a slightly different but equivalent question of “how to scale A
805 such that Av has entry size $\Theta(1)$?”

806 G.1.1 Preparation for the Derivations

807 By the results of [40], each (pre-)activation vector and its gradient vector in a multi-layer perceptron
808 have approximately iid coordinates in the large width limit,²⁶ and something similar can be said for
809 more advanced networks such as ResNet and Transformers²⁷. In particular, to each such vector v ,
810 we can associate a random variable Z^v that represents the coordinate distribution of v . If vector u is
811 correlated with v , then Z^u will also be correlated with Z^v , and $\lim_{n \rightarrow \infty} v^\top u/n = \mathbb{E} Z^u Z^v$.

²⁴in the sense that the the variance of the entries are $\Theta(1)$

²⁵Here “correlated” formally means v depends on W^\top in a Tensor Program. This essentially captures all scenarios of “ v correlated with W ” that occurs in deep learning.

²⁶Our intuition here is derived from the assumption that width is much larger than training time; of course, as illustrated by our myriad experiments, these intuition are very useful even when this is not the case, such as when training to convergence.

²⁷E.g. in a convnet, the (pre-)activations are iid across channels, but correlated across pixels

812 **G.1.2 Linear Tensor Product Matrix (e.g. SGD Updates)**

813 The case of (linear) tensor product matrix can be reduced to the outer product case by linearity. Given
 814 $u, v, x \in \mathbb{R}^n$ having approximately iid coordinates (of size $\Theta(1)$) like so, we can form the outer
 815 product

$$A \stackrel{\text{def}}{=} u \otimes v/n = uv^\top/n, \quad (5)$$

816 which is the form of a single (batch size 1) gradient update to a weight matrix. Then, by Law of
 817 Large Numbers,

$$Ax = u \frac{v^\top x}{n} \approx cu, \quad \text{where } c = \mathbb{E} Z^v Z^x.$$

818 So Ax also has approximately iid coordinates, distributed like $Z^{Ax} \stackrel{\text{def}}{=} Z^u \mathbb{E} Z^v Z^x$. Likewise, if A is
 819 a sum of outer products $A = \sum_{i=1}^k u^i \otimes v^i/n$, then

$$Ax = \sum_{i=1}^k u^i \frac{v^{i\top} x}{n}, \quad \text{with coordinates distributed as } Z^{Ax} = \sum_{i=1}^k Z^{u^i} \mathbb{E} Z^{v^i} Z^x.$$

820 Notice that each coordinate of A has size $\Theta(1/n)$. The above reasoning shows that, in order for Ax
 821 to have coordinate size $\Theta(1)$ (assuming x does), then $\Theta(1/n)$ is the right coordinate size for A , in
 822 the general case that v^i and x are correlated (as is generically the case during gradient descent, with
 823 $A = \Delta W$ for some weights W and x being the previous activations).²⁸

824 **G.1.3 Nonlinear Tensor Product Matrix (e.g. Adam Updates)**

825 When using Adam or another adaptive optimizer that normalizes the gradient coordinatewise before
 826 applying them, we need to modify our argument slightly to obtain the right coordinate size scaling of
 827 the matrix. The gradient update A , after such normalization, will take the form of

$$A_{\alpha\beta} = \psi(u_\alpha^1, \dots, u_\alpha^k, v_\beta^1, \dots, v_\beta^k), \quad \text{for some } \psi : \mathbb{R}^{2k} \rightarrow \mathbb{R} \text{ and vectors } u^i, v^j \in \mathbb{R}^n. \quad (6)$$

828 We say a matrix of this form is a *nonlinear tensor product matrix*.

829 First, note the tensor product matrices (e.g. the form of SGD update) discussed previously (Eq. (5))
 830 already takes this form, with $\psi(u_\alpha^1, \dots, u_\alpha^k, v_\beta^1, \dots, v_\beta^k) = n^{-1}(u_\alpha^1 v_\beta^1 + \dots + u_\alpha^k v_\beta^k)$, so Eq. (6)
 831 is a strict generalization of linear tensor products. Next, for the example of Adam, each gradient
 832 update is μ/σ where μ (resp. σ^2) is the moving average of previous (unnormalized) gradients (resp.
 833 the coordinatewise square of the same).²⁹ If these unnormalized gradients are the outer products
 834 $u^1 \otimes v^1, \dots, u^k \otimes v^k$, then the update has coordinates

$$(\mu/\sigma)_{\alpha\beta} = \psi(u_\alpha^1, \dots, u_\alpha^k, v_\beta^1, \dots, v_\beta^k) \stackrel{\text{def}}{=} \sum_i \gamma_i u_\alpha^i v_\beta^i / \sqrt{\sum_i \omega_i (u_\alpha^i v_\beta^i)^2}, \quad (7)$$

835 where γ_i and ω_i are the weights involved in the moving averages.

836 Now suppose we have some $A \in \mathbb{R}^{n \times n}$ of the form Eq. (6), where $u^i, v^i \in \mathbb{R}^n$ have approximately
 837 iid coordinates (of size $\Theta(1)$), and $\psi = n^{-1}\bar{\psi}$ where $\bar{\psi}$ doesn't depend on n (in terms of Adam where
 838 $\bar{\psi}$ corresponds to the ψ of Eq. (7), this corresponds to using a learning rate of $1/n$). Then for $x \in \mathbb{R}^n$
 839 having approximately iid coordinates of size $\Theta(1)$, by Law of Large Numbers,

$$(Ax)_\alpha = \frac{1}{n} \sum_{\beta=1}^n \bar{\psi}(u_\alpha^1, \dots, u_\alpha^k, v_\beta^1, \dots, v_\beta^k) x_\beta \approx \mathbb{E} \bar{\psi}(u_\alpha^1, \dots, u_\alpha^k, Z^{v^1}, \dots, Z^{v^k}) Z^x \stackrel{\text{def}}{=} \Psi(u_\alpha^1, \dots, u_\alpha^k).$$

840 Here we made the obvious definition

$$\Psi : \mathbb{R}^k \rightarrow \mathbb{R}, \quad \Psi(r_1, \dots, r_k) \stackrel{\text{def}}{=} \mathbb{E} \bar{\psi}(r_1, \dots, r_k, Z^{v^1}, \dots, Z^{v^k}) Z^x.$$

²⁸In some corner cases when x is uncorrelated with v , then $v^\top x = \Theta(\sqrt{n})$ by Central Limit, so actually Ax has $\Theta(1/\sqrt{n})$ coordinates. However, this case does not come up much in the context of training neural networks.

²⁹Adam also has bias correction for the moving averages which can be accommodated easily, but for simplicity we omit them here.

841 Thus Ax also has approximately iid coordinates (of size $\Theta(1)$),

$$Z^{Ax} \stackrel{\text{def}}{=} \Psi(Z^{u^1}, \dots, Z^{u^k}).$$

842 For example, in the SGD example with $A = u \otimes v/n$ and $\bar{\psi}(u_\alpha, v_\beta) = u_\alpha v_\beta$, this formula gives
843 $Z^{Ax} = \Psi(Z^u)$ where $\Psi(z) = z \mathbb{E} Z^v Z^x$, recovering the earlier derivation.

844 In any case, the point here is that A has coordinate size $\Theta(1/n)$, and this is the unique scaling that
845 leads to Ax having coordinate size $\Theta(1)$.

846 **G.1.4 Vector Case (e.g. Readout Layer)**

847 The vector A case is similar to the tensor product cases above.

848 **G.1.5 Gaussian Matrix (e.g. Hidden Weights Initialization)**

849 Now consider the case where $A \in \mathbb{R}^{n \times n}$ is random Gaussian matrix with $A_{\alpha\beta} \sim \mathcal{N}(0, 1/n)$ and
850 $x \in \mathbb{R}^n$ has approximately iid coordinates distributed like Z^x . In the context of neural network
851 training, A should be thought of as a randomly initialized weight matrix, and x for example can be
852 taken to be an activation vector in the first forward pass.

853 If x is independent from A (or sufficiently uncorrelated), then each coordinate $(Ax)_\alpha$ has variance
854 $\mathbb{E}(Z^x)^2 = \Theta(1)$ (so by definition has size $\Theta(1)$). Thus, here A having $\Theta(1/\sqrt{n})$ coordinates leads
855 to Ax having $\Theta(1)$ coordinates, in contrast to the tensor product case above.

856 When x is correlated with A , it turns out the same scaling applies ($\Theta(1/\sqrt{n})$ is the unique scaling for
857 A 's entries such so that Ax has $\Theta(1)$ entries), but the reasoning is much more subtle: In the context
858 of neural network training, it turns out all scenario where x is correlated with A can be reduced
859 to the case where $x = \phi(A^\top y, \dots)$ for some coordinatewise nonlinearity ϕ and some other vector
860 \mathbb{R}^n .³⁰ Let's consider a very simple example with $x = A^\top \mathbf{1}$ for the all 1s vector $\mathbf{1} \in \mathbb{R}^n$ (which has
861 coordinate size $\Theta(1)$ as can be checked easily). Then, for each index $\alpha \in [n]$, we can calculate

$$(AA^\top \mathbf{1})_\alpha = \sum_{\beta, \gamma} A_{\alpha\beta} A_{\gamma\beta} = \sum_{\beta} A_{\alpha\beta}^2 + \sum_{\beta} \sum_{\gamma \neq \alpha} A_{\alpha\beta} A_{\gamma\beta}.$$

862 Since $\mathbb{E} A_{\alpha\beta}^2 = 1/n$, by the Law of Large Number, the first sum $\sum_{\beta} A_{\alpha\beta}^2 \approx 1$. On the other hand,
863 there are n summands of the form $\sum_{\gamma \neq \alpha} A_{\alpha\beta} A_{\gamma\beta}$, all iid with variance $\frac{n-1}{n^2} = \Theta(1/n)$. Thus by
864 the Central Limit Theorem, we expect $\sum_{\beta} \sum_{\gamma \neq \alpha} A_{\alpha\beta} A_{\gamma\beta} \approx \mathcal{N}(0, 1)$. Therefore, each coordinate
865 of $(AA^\top \mathbf{1})_\alpha$ looks like $1 + \mathcal{N}(0, 1) = \mathcal{N}(1, 1)$ and thus has size $\Theta(1)$; again this is caused by A
866 having $\Theta(1/\sqrt{n})$ coordinates.

867 This example can be generalized to more general x that is correlated with A , but the mathematics is
868 quite involved. See [39] for more details.

869 **G.2 Deriving μP for Any Architecture**

870 Armed with the insight from the last section, we now outline the key steps to derive μP for any
871 architecture. In practice, μP of [40] implies the following desiderata

872 **Desiderata G.1.** At any time during training

- 873 1. Every (pre)activation vector in a network should have $\Theta(1)$ -sized coordinates³¹
- 874 2. Neural network output should also be $\Theta(1)$.
- 875 3. All parameters should be updated as much as possible (in terms of scaling in width) without
876 leading to divergence.

877 Let's briefly justify these desiderata. For the desideratum 1, if the coordinates are $\omega(1)$ or $o(1)$,
878 then for sufficiently wide networks their values will go out of floating point range. This problem is

³⁰This is because every "reasonable" deep learning computation can be expressed in a Tensor Program.

³¹In a convnet, a (pre-)activation vector corresponds to a single pixel across all channels; in general, we expect (pre-)activations are iid across channels, but correlated across pixels

879 particularly acute for low-precision formats that are essential for training large models such as BERT
 880 or GPT. Moreover, a general nonlinearity is only well-behaved if its input is in a fixed range (although
 881 this is not a problem for homogeneous nonlinearities like relu). For example, for tanh nonlinearity, if
 882 the preactivation is vanishing $o(1)$, then tanh is essentially linear; if the preactivation is exploding
 883 $\omega(1)$, then the tanh gradient vanishes.

884 For the desideratum 2, a similar justification applies to the numerical fidelity of the loss function and
 885 loss derivative.

886 Finally, desideratum 3 means that 1) we are doing “maximal feature learning” [40] and 2) every
 887 parameter contribute meaningfully in the infinite-width limit. This ensures that learning rate “plays
 888 the same role” in the finite-width case as in the infinite-width limit. For example, it prevents the
 889 scenario where a weight matrix gets stuck at initialization in the limit for any learning rate (so
 890 learning rate does not matter) but evolves nontrivially in any finite-width network (so learning rate
 891 does matter).

892 These desiderata will essentially uniquely single out μP . More formally, μP is the unique parametriza-
 893 tion that admits feature learning in all parameters of the neural network [40], and this property
 894 theoretically guarantees hyperparameter transfer across width (for sufficiently large width). However,
 895 for the sake of reaching a broader audience, we will focus more on the intuitive derivations from the
 896 desiderata rather than on this formal aspect.

897 Below, we assume for simplicity that the width of every layer is n , and we focus only on dense
 898 weights. Later, we will discuss convolutions and varying the widths between layers.

899 G.2.1 μP Derivation From the Desiderata

900 **Output Weights** Suppose $W \in \mathbb{R}^{1 \times n}$ is an output weight. By desideratum 1, the input x to W
 901 has $\Theta(1)$ -sized coordinates. Thus W should have $\Theta(1/n)$ coordinates so that $Wx = \Theta(1)$. We
 902 can initialize W with $\Theta(1/n)$ coordinates and scale its (per-layer) LR so that ΔW has $\Theta(1/n)$
 903 coordinates as well. However, in order to use the same SGD learning rate for all layers, we instead
 904 equivalently 1) reparametrize $W = \frac{1}{\sqrt{n}}w$ with trainable w , 2) initialize w with $\Theta(1/\sqrt{n})$ coordinates,
 905 and 3) use $\Theta(1)$ learning rate (of w) for SGD. For Adam, however, we cannot use the same learning
 906 rate for every layer, so we set the Adam learning rate of w to be $\Theta(1/\sqrt{n})$.

907 **Hidden Weights** Consider a square weight matrix $W \in \mathbb{R}^{n \times n}$. Desiderata 1 guarantees that the
 908 input x to W has $\Theta(1)$ -sized coordinates. Generally, x will be correlated with W . By Table 9, we
 909 can immediately derive

910 **Initialization** W should be randomly initialized with coordinate size $\Theta(1/\sqrt{n})$

911 **LR** The learning rate should be scaled so that ΔW has coordinate size $\Theta(1/n)$

912 so that $(W_0 + \Delta W)x$ is $\Theta(1)$ if x is, inductively satisfying desideratum 1. With Adam, this just
 913 means the per-layer LR is $\Theta(1/n)$. With SGD and the scaling of output layers above, we can calculate
 914 that the gradient of W has $\Theta(1/n)$ coordinates, so the $\Theta(1)$ SGD LR derived above suffices as well.

915 **Input Weights** Suppose $W \in \mathbb{R}^{n \times d}$ is an input weight. To satisfy desiderata 1 (i.e. for any input ξ ,
 916 $W\xi$ should have $\Theta(1)$ coordinates), we want W to have $\Theta(1)$ coordinates. We can initialize W with
 917 $\Theta(1)$ coordinates and scale its (per-layer) LR so that ΔW has $\Theta(1)$ coordinates as well. However,
 918 in order to use the same SGD learning rate for all layers, we instead 1) reparametrize $W = \sqrt{n}w$
 919 with trainable w , 2) initialize w with $\Theta(1/\sqrt{n})$ coordinates, and 3) use $\Theta(1)$ for SGD learning rate.
 920 Again, for Adam LR, we have to use a per-layer learning rate (of w), for which we take $\Theta(1/\sqrt{n})$.

921 **Biases** Biases follow the same reasoning as input weights (just think of it as an input weight with
 922 input 1).

923 **Attention** Suppose the key dimension d_k is tending to infinity with width with number of heads
 924 n_{head} fixed. Then the key-query contraction $q^\top k \in \mathbb{R}$ scales like $\Theta(d_k)$ by Law of Large Numbers
 925 (instead of Central Limit Theorem because q and k are generally correlated) and desideratum 1, hence
 926 the $1/d_k$ we propose rather than $1/\sqrt{d_k}$.

927 Now suppose instead that n_{head} tends to infinity with width with d_k fixed. Let $K, Q \in$
 928 $\mathbb{R}^{N \times d_k \times n_{head}}$, $V \in \mathbb{R}^{N \times d_v \times n_{head}}$ be keys, queries, and values across all heads and tokens. Think-

929 ing of $N \times d_k$ as constants, we may view attention as a nonlinearity coordinatewise in the n_{head}
930 dimension. Then it’s clear that our parametrization described above already works.

931 Finally, we may freely let d_k and n_{head} both tend to infinity, and the above reasoning shows that our
932 parametrization still works.

933 **Changing Width Ratios** As noted above, at any time in training, every (pre-)activation vector will
934 have approximately iid coordinates (of order $\Theta(1)$ by desideratum 1). Another desideratum for μP is
935 to ensure that this coordinate distribution (at any particular time) stays roughly invariant as widths
936 increases. When all layer widths are tied, this is automatic if the other desiderata are satisfied, hence
937 why we did not list this above.

938 When width ratios vary, this is not automatic. In this case, we need to choose whether to replace each
939 n with fan-in or fan-out (or some function of them). Making the wrong choices will let the coordinate
940 distributions vary with width ratios.

941 It turns out the correct choice, as shown in Table 3, is to replace n with fan-in for the input layers
942 and with fan-out for the output layers. For the hidden weights, we replace n with fan-in so that
943 the forward pass is preserved. When using Adam (and assuming the initialization of W is quickly
944 dominated by the change in W), this ensures that the (pre-)activation coordinate distributions are
945 preserved at any time during training even if we vary widths in different layers differently. (For
946 SGD this doesn’t quite work in general because the varying width ratios change the gradient sizes of
947 different layers differently, whereas Adam always normalizes the gradient coordinatewise).

948 **Convolution** A convolution weight tensor $W \in \mathbb{R}^{\text{fan_out} \times \text{fan_in} \times s_1 \times s_2}$ with kernel size $s_1 \times s_2$
949 can be thought of just as a $s_1 s_2 = \Theta(1)$ -sized collection of $\text{fan_out} \times \text{fan_in}$ dense weights. Then
950 all of our discussions above apply accordingly.

951 G.3 Why Other Parametrizations Cannot Admit Hyperparameter Transfer

952 **Standard Parametrization (SP)** SP doesn’t work essentially because it leads to blow-up in the
953 infinite-width limit.

- 954 1. For Adam (with LR $\Theta(1)$), ΔW would have $\Theta(1)$ coordinates, causing preactivations to
955 blow up like $\Theta(n)$ by Desideratum 1 and Table 9.
- 956 2. For SGD, the gradient of $\mathbb{R}^{n \times n}$ weight has $\Theta(1/\sqrt{n})$ coordinates because the last layer
957 has no $1/\sqrt{n}$ factor, so $\Theta(1)$ learning rate would make preactivation scale like $\Theta(\sqrt{n})$ and
958 hence blow up.

959 If we use $\Theta(1/\text{width})$ learning rate, then blow-up does not occur. However, this infinite-width limit
960 is in the kernel regime and thus does not allow hyperparameter transfer for the same reason that NTP
961 below does not.

962 **Neural Tangent Parametrization (NTP)** We have concrete examples, e.g. Word2Vec in [40],
963 where the NTK limit has trivial performance — so hyperparameters have no effect at all — vastly
964 outperformed by finite-width networks — where hyperparameters matter. More importantly, wider
965 does not always do better in NTP, especially in tasks where feature learning is crucial [40]. So in the
966 context of modern deep learning e.g. large language model pretraining, NTP (or SP with $\Theta(1/\text{width})$
967 LR) does not make sense for wide neural networks.

968 **Other Parametrizations** Recall the *Dynamical Dichotomy Theorem* proven in [40], which says
969 that any nontrivial stable “natural parametrization” (formally, “*abc-parametrization*,” [40]) either
970 admits a feature learning limit or a kernel limit, but not both.

971 Our argument above against SP and NTP will also work against any parametrization inducing a
972 kernel limit. Therefore, it remains to ask, can other *feature learning* parametrizations transfer
973 hyperparameters?

974 We argue no. As shown in [40], any other feature learning parametrization differs from μP essentially
975 only in that some parameters are not updated maximally. By [40, Sec 6.4], in the infinite-width limit,
976 such parameters can be thought of as being fixed at initialization. Therefore, in such infinite-width
977 limits, the learning rate of such parameters becomes useless. Therefore, we cannot hope for the

978 hyperparameter landscape of the limit to reflect the hyperparameter landscape of finite-width neural
979 networks.

980 μP is the unique feature learning parametrization that updates all parameters maximally, so that the
981 learning rate of each parameter plays approximately the same role in finite-width neural networks as
982 in the infinite-width limit. Consequently, the hyperparameter landscape of the μP limit should reflect
983 the hyperparameter landscape of finite-width neural networks.

984 **H Nuances of the Hyperparameter Landscape Convergence Intuition**

985 **What converges and what doesn't** We want to tune hyperparameters on a small model with width
986 N such that its hyperparameter landscape looks like that of a large model with width $\gg N$. However,
987 for this to be useful, we *do not want* the small model (as a function) after training to be close to that
988 of the large model — otherwise there is no point in training the large model to begin with. So N 1)
989 must be large enough so that the hyperparameter optimum converges, but 2) cannot be so large that
990 the functional dynamics converges. The fact that such N exists, as demonstrated by our experiments,
991 shows that: In some sense, the hyperparameter optimum is a “macroscopic” or “coarse” variable
992 which converges quickly with width, while the neural network function is a very “microscopic” or
993 “fine” detail that converges much more slowly with width.

994 The same discussion applies to scaling of batch size as well [21].