

Tensor Programs V: Tuning Large Neural Networks via Zero-Shot Hyperparameter Transfer

Greg Yang

in Collaboration with Edward Hu, Igor Babuschkin, Szymon Sidor, David Farhi, Nick Ryder,

Jakub Pachocki, Xiaodong Liu, Weizhu Chen, Jianfeng Gao



Explore with a small model



Model fails to train when scaled up with the same hyperparameters



An Example of What *This Work* Allows You To Do



Before

After

Enter Maximal Update Parametrization, abbreviated μP

A principled way of scaling initialization and learning rate with width

	Input weights & all biases	Output weights	Hidden weights
Init. Var.	$\begin{array}{c} 1/_{\text{fan_in}}\\ \text{fan_out} & (1)\\ 1 \end{array}$	$1/fan_in^2$ (1/fan_in)	¹ /fan_in
SGD LR		1/fan_in (1)	1
Adam LR		1/fan_in (1)	1/fan_in (1)

Note: focus on scaling with fan_in or fan_out; everything else is a tunable constant

Key Property of μ P: Hyperparameter Stability





$\mu {\rm Transfer:} \ {\rm Zero-Shot} \ {\rm Hyperparameter} \ {\rm Transfer}$



"Transfer" = optimal hyperparameter remains stable with model size

CAREFUL VALIDATION ON MODELS OF INCREASING SCALE



PRACTICAL BENEFITS OF μ **P**

Preserves hyperparameter optimum across width

- Allows zero-shot hyperparameter transfer
 - Efficient tuning
 - Can tune enormous models only on a single GPU
 - Very fast
 - The only method for tuning large models
- Altogether enable *more reliable scaling-up* of neural networks

WHAT DO THESE HAVE IN COMMON?



- Impressive achievements of humankind
- Each empirical test is very expensive
- Require extensive theoretical calculation first before launching any empirical test

ECONOMY OF INSIGHTS AT SCALE





What is Maximal Update Parametrization?

 $n = width, d_{in} = data dimension$

WARMUP: MLP

Standard Parametrization (SP)

$$\begin{split} f(\xi) &= W^{3\top} \phi(W^{2\top} \phi(W^{1\top} \xi + b^1) + b^2) \\ \text{with init.} \quad W^1 \sim \mathcal{N}(0, 1/d_{in}), \ W^{\{2,3\}} \sim \mathcal{N}(0, 1/n), \ b^{\{1,2\}} = 0, \quad \mathsf{LR} \ \eta \end{split}$$

 μ -Parametrization (μ P)

initialize $W^1 \sim \mathcal{N}(0, 1/d_{in}), W^2 \sim \mathcal{N}(0, 1/n), W^3 \sim \mathcal{N}(0, 1/n^2), b^{\{1,2\}} = 0$ with SGD learning rates $\eta_{W^1} = \eta_{b^1} = \eta_{b^2} = \eta n, \ \eta_{W^2} = \eta, \ \eta_{W^3} = \eta n^{-1}.$



 $n = width, d_{in} = data dimension$

WARMUP: MLP

 μ -Parametrization (μ P)



- didn't think about correlation during training!



$\mu\text{-}\mathsf{PARAMETRIZATION}$ IN GENERAL

	Input weights & all biases	Output weights	Hidden weights
Init. Var. SGD LR Adam LR	$\begin{array}{c} 1/_{\text{fan_in}}\\ \text{fan_out} & (1)\\ 1 \end{array}$	$1/_{fan_in^2}$ (1/_{fan_in}) 1/_{fan_in} (1) 1/_{fan_in} (1)	$\frac{1/\text{fan_in}}{1}$ $\frac{1}{1/\text{fan_in}}$ (1)

Note: focus on scaling with fan_in or fan_out; everything else is a tunable constant

For Transformers, also use 1/d instead of $1/\sqrt{d}$ attention, i.e. the attention logit is calculated as $k^{T}q/d$ instead of $k^{T}q/\sqrt{d}$ where query q and key k have dimension d.

Which Hyperparameters Can Be μTransferred? 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*

- Why not regularization?
 - The amount of regularization (for the purpose of controlling overfitting) naturally depends on both the model size and data size
 - so we should not expect transfer to work if the parametrization only depends on model size (width)
 - Therefore, if regularization is the bottleneck of performance of a large model, $\mu {\rm P}$ may not help with that
 - Luckily, for large scale pretraining, we still have too much data so regularization is not a bottleneck
- While transfer across width is justified by theory, we empirically explore transferring across other scale dimensions like depth, batch size, training time, seq length.
 - But without theoretical understanding, we cannot be as confident that these results hold beyond the *reasonable* scales and the problem settings we did experiments in

Verifying Transfer Across Width



What if We Change Width Ratio?



Figure 12: Learning rate landscape in μ P is stable even if we vary d_{ffn} by a factor of 32, fixing d_{model} .

You Can Vary Both d_{head} and n_{head}





Figure 13: μ 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.

WIDER IS BETTER THROUGHOUT TRAINING



Figure 9: Wider is always better in training loss under μ P, but not in SP, given the same HP. Learning curves for μ P and SP with different learning rates, aggregated over 5 seeds. (Left) Wider μ P models always achieve better training loss at any time in training. (Middle) If using a small learning rate, SP models can appear to do so up to some large width, at which point the pattern fails (at width 2048 in our plot). (Right) If using a large learning rate, SP model can do *worse* with width; here the SP model is identical to the μ P model in (Left) at width 128.

EXPLORING EMPIRICAL TRANSFER ACROSS DEPTH









EXPLORING EMPIRICAL TRANSFER ACROSS OTHER DIMENSIONS



preLN Transformer

EXPLORING EMPIRICAL TRANSFER ACROSS OTHER DIMENSIONS



postLN Transformer

Empirical Efficacy of μ **Transfer**

CAREFUL VALIDATION ON MODELS OF INCREASING SCALE



Transformer on IWSLT14 De-En



 μ P + μ Transfer



Standard Parametrization

			Val. BLEU Percentiles			
Setup	Total Compute	#Samples	25	50	75	100
fairseq 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	
Transfer from 0.25x	1x	64	35.27	35.33	35.45	35.53

log₂LearningRate

Performance-Compute Pareto Frontier



IWSLT Transformer is still fairly small --- efficiency only goes up with model scale!

CAREFUL VALIDATION ON MODELS OF INCREASING SCALE



CAREFUL VALIDATION ON MODELS OF INCREASING SCALE



Tuning BERT with μ Transfer

<u>Step 1:</u>

Parameterize BERT in μP

<u>Step 2:</u>

Tune hyperparameters on BERT_{SMALL} via random search (256 combinations), training only for 10% of tokens

<u>Step 3:</u>

Copy the best hyperparameter combination to $\mathsf{BERT}_{\mathsf{BASE}}$ and $\mathsf{BERT}_{\mathsf{LARGE}}$

- \checkmark Tune once, use for a family of models
- \checkmark Only run the large models once

The total tuning cost is equivalent to 1 full training run of BERT-large

Model	Method	Model Speedup	Total Speedup	Test loss	MNLI (m/mm)	QQP
$egin{array}{c} {\sf BERT}_{base} \ {\sf BERT}_{base} \ {\sf BERT}_{base} \end{array}$	Megatron Default Naive Transfer µTransfer (Ours)	1x 4x 4x	1x 40x 40x	1.995 tı 1.970	84.2/84.2 raining diverged 84.3/84.8	90.6 90.8
$egin{array}{l} {\sf BERT}_{large} \ {\sf BERT}_{large} \ {\sf BERT}_{large} \end{array}$	Megatron Default Naive Transfer µTransfer (Ours)	1x 22x 22x	1x 220x 220x	1.731 t. 1.683	86.3/86.2 raining diverged 87.0/86.5	90.9 91.4

CAREFUL VALIDATION ON MODELS OF INCREASING SCALE



CAREFUL VALIDATION ON MODELS OF INCREASING SCALE



OpenAl GPT-3 Family + μ **P**

Hyperparameter Optimum is Stable



Wider is Always Better Given the Same HPs



OpenAl GPT-3 6.7B + μ **Transfer**



Total tuning compute budget is only 7% of training budget!!!

Theoretical Foundation

WHAT IS A PARAMETRIZATION?

	Input weights & all biases	Output weights	Hidden weights
Init. Var.	?	?	?
SGD LR	?	?	?
Adam LR	?	?	?

- Tells you how to scale HP when model size changes
 - "half learning rate when width doubles"
- *Doesn't* tell you how to set specific values of HP at particular model sizes
 - "use learning rate 1e-3 at width 1024"
- For any fixed model size, all parametrizations are equivalent up to tunable HP
 - Parametrization gives a way for expressing hyperparameters, just like a basis gives a way to express vectors in a vector space
 - A priori, there is no canonical parametrization, just like a priori there is no canonical basis in a vector space

EACH PARAMETRIZATION CORRESPONDS TO A ∞ -WIDTH LIMIT



OPTIMAL SCALING THESIS

Verify empirically

Verify theoretically

HP Transfer

A parametrization preserves HP optimum across model sizes

Maximal Limit

Its limit maximizes feature learning (without blowing up)

This is a thesis for any notion of model size: width, depth, #experts, etc

OPTIMAL SCALING THESIS



"PROVING" OPTIMAL SCALING THESIS FOR WIDTH



Dynamical Dichotomy Theorem

A Caricature of Space of Parametrizations



٠

- Feature learning
- Function evolution cannot be described purely in the function space

"PROOF" OF OPTIMAL SCALING THESIS FOR WIDTH

- Obviously unstable or trivial parametrizations won't work
 - E.g., optimal learning rate will converge to 0 or diverge.
- Kernel regime parametrizations have trivial performance on pretraining tasks
 - E.g., NTK performs trivially on Word2vec; see TP4
- Nonmaximal feature learning limits differ from μP's (maximal) limit only in that some parameters are stuck at init, ignoring the learning rate
 - But in finite networks, for such parameters, learning rate obviously has nontrivial effect --- so there's no way we can expect transfer
- The only parametrization left is μP

A Caricature of Space of Parametrizations



TENSOR PROGRAMS



WHAT IS A TENSOR PROGRAM?

• A computation composed of 1) matrix multiplication and 2) coordinatewise nonlinearities



Any "useful" deep learning computation can be expressed as a Tensor Program.



Autograd	gradients
----------	-----------

Master Theorem +	→ ∞-width limit
------------------	-----------------





CONCLUSION



LOOKING FORWARD: THE THEORY OF EVERYTHING

- μP solves the transfer problem for width in a principled way. Can we do it for all other compute hyperparameters?
 - Ex: depth, #experts, training time, sequence length
 - Naïve transfer seems to work OK empirically, but as we go to larger scales, likely they will break down
 - Verify the Optimal Scaling Thesis
- This is akin to the search for the Theory of Everything in physics
 - The different scale hyperparameters of a neural network are like the different fundamental forces of nature
 - We look for a theory that tells us how to scale every aspect of a model together, just like physicists look for a GUT unifying all fundamental forces of nature



Pytorch library for μP

> pip install mup

OR

> git clone github.com/microsoft/mup