

μ LO: Compute-Efficient Meta-Generalization of Learned Optimizers

Benjamin Thérien^{1,2}

BENJAMIN.THERIEN@MILA.QUEBEC

Charles-Étienne Joseph^{1,2}

CHARLES-ETIENNE.JOSEPH@MILA.QUEBEC

Boris Knyazev⁴

KNYAZEVB@MILA.QUEBEC

Edouard Oyallon⁵

EDOUARD.OYALLON@CNRS.FR

Irina Rish^{1,2}

IRINA.RISH@MILA.QUEBEC

Eugene Belilovsky^{2,3}

EUGENE.BELILOVSKY@CONCORDIA.CA

¹Université de Montréal; ²Mila – Quebec AI Institute; ³Concordia University, Montréal;

⁴Samsung - SAIT AI Lab, Montréal; ⁵Flatiron Institute, New York, NY, USA.

Abstract

Learned optimizers (LOs) can significantly reduce the wall-clock training time of neural networks, substantially reducing training costs. However, they can struggle to optimize unseen tasks (meta-generalize), especially when training networks much larger than those seen during meta-training. To address this, we derive the Maximal Update Parametrization (μ P) for two popular learned optimizer architectures and propose a simple meta-training recipe for μ -parameterized LOs (μ LOs). Our empirical evaluation demonstrates that LOs meta-trained with our recipe substantially improve meta-generalization to wider unseen tasks when compared to LOs trained under standard parametrization (e.g., as they are trained in existing work). When applying our μ LOs, each trained for less than 250 GPU-hours, to large-width models we are often able to match or exceed the performance of pre-trained VeLO, the most performant publicly available learned optimizer, meta-trained with 4000 TPU-months of compute. We also empirically observe that learned optimizers trained with our μ LO recipe also exhibit substantially improved meta-generalization to deeper networks ($5\times$ meta-training) and remarkable generalization to much longer training horizons ($25\times$ meta-training).

1. Introduction

Deep learning’s success can, in part, be attributed to its ability to learn effective representations for downstream tasks. Notably, this resulted in the abandonment of a number of heuristics (e.g., hand-designed features in computer vision [8, 17]) in favor of end-to-end learned features. However, one aspect of the modern deep-learning pipeline remains hand-designed: gradient-based optimizers. While popular optimizers such as Adam or SGD provably converge to a local minimum in non-convex settings [13, 15, 24], there is no reason to expect these hand-designed optimizers reach the global optimum at the optimal rate for a given problem. Given the lack of guaranteed optimality and the clear strength of data-driven methods, it is natural to turn towards data-driven solutions for improving the optimization of neural networks.

To improve hand-designed optimizers, [3, 18, 19, 30] replaced them with small neural networks called learned optimizers (LOs). [20] showed that scaling up the training of such optimizers can significantly improve wall-clock training speeds and supersede existing hand-designed optimizers. However, LOs have limitations in *meta-generalization* – optimizing new problems. For example, despite training for 4000 TPU months, VeLO [20] is known to (1) generalize poorly to longer optimization problems (e.g., more steps) than those seen during meta-training and (2) have difficulty

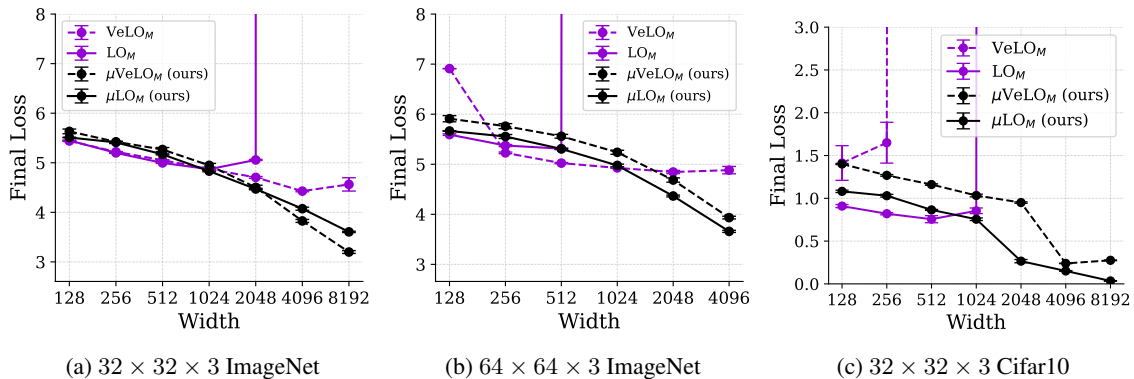


Figure 1: **Generalization beyond meta-training widths is severely limited without our approach.** We report the final loss after 1000 steps (e.g., the inner problem length used when meta-training) for models of different widths. Each point is the average final training loss over 5 seeds with standard error bars. We observe that both μ LOs consistently obtain lower loss values as the tasks become wider. In contrast, their SP LO counterparts either diverge before reaching 1000 steps on the wider tasks or make little progress as width increases.

optimizing models much larger than those seen during meta-training. Given the high cost of meta-training LOs (e.g., when meta-training, a *single training example* is analogous to training a neural network for many steps), it is essential to be able to train learned optimizers on small tasks and generalize to larger ones. [11] explore preconditioning methods to improve the generalization from shorter to longer optimization problems (e.g., ones with more steps). However, no works have tackled the meta-generalization of LOs to wider models in a principled way.

To address the meta-generalization problem of LOs, we recognize that this problem can be reformulated as *zero-shot hyperparameter transfer* [36]. The latter involves selecting optimal hyperparameters of hand-designed optimizers for training very large networks (that one cannot afford to tune directly) by transferring those tuned on smaller versions of the model. Under the standard parametrization (SP)¹, the optimal hyperparameters of an optimizer used for a small model do not generalize well to larger versions of the model. However, when a small model is tuned using the Maximal Update Parametrization (μ P), and its larger counterparts are also initialized with μ P, the small and large models share optimal hyperparameters [36]. Given the appealing connection between zero-shot hyperparameter transfer in hand-crafted optimizers and meta-generalization in LOs, we ask the following questions: *Can learned optimizers be meta-trained under μ P? How would the resulting optimizers perform on wider unseen tasks?* We seek to answer these questions in the following study. Specifically, we consider the recent LO architectures [19, 20] and demonstrate that μ P can be adapted to these optimizers leading to our μ LO optimizers. We subsequently conduct an empirical evaluation that reveals the power of our μ LOs and their advantages for scaling learned optimizers.

Our contributions can be summarized as follows:

- We derive μ -parameterization for two popular learned optimizer architectures (VeLO and small_fc_lopt) and propose a training recipe for μ LOs.

1. When we refer to SP, we follow the same meaning as [36]. That is, we consider SP to designate a parameterization that does not admit HP transfer. However, we note that recent work [10] shows hyperparameter transfer is possible in SP.

- We demonstrate that μ LOs meta-trained with our recipe significantly improve generalization to wider networks when compared to their SP counterparts and several strong baselines and that, for wider counterparts of the meta-training tasks, they outperform VeLO (meta-trained with 4000 TPU-months of compute).
- We demonstrate empirically that μ LOs meta-trained with our recipe show improved generalization to deeper networks ($5\times$ meta-training) when compared to their SP counterparts.
- We demonstrate empirically that μ LOs meta-trained with our recipe significantly improve generalization to longer training horizons ($25\times$ meta-training) when compared to their SP counterparts.

Our results show that μ LOs significantly improve learned optimizer generalization without increasing meta-training costs. This constitutes a noteworthy improvement in the scalability of meta-training LOs.

2. μ -parametrization for Learned Optimizers

Parameterizing an optimizee neural network in μ P requires special handling of the initialization variance, pre-activation multipliers, and optimizer update for each weight matrix $w \in \mathbb{R}^{n \times m}$ in the network. Specifically, these quantities will depend on the functional form of the optimizer and the dependence of n (FAN_OUT) and m (FAN_IN) on width. We will refer to weight matrices in a network of width h as hidden layers if $\Theta(n) = \Theta(m) = \Theta(h)$, as output layers if $\Theta(n) = 1, \Theta(m) = \Theta(h)$, and as input layers if $\Theta(n) = \Theta(h), \Theta(m) = \Theta(h)$.

Consider a model *to be optimized* g_w with weights in layers l denoted w_l . We apply and construct μ LOs as follows.

Initialization- μ . w_l which are hidden and input layers have their weights initialized as $\mathcal{N}(0, \frac{1}{\sqrt{\text{FAN_IN}}})$. While output layers have their weights initialized as $\mathcal{N}(0, 1)$.

Multipliers- μ . Output layer pre-activations are multiplied by $\frac{1}{\text{FAN_IN}}$ during the forward pass.

Updates- μ . The update by f_ϕ on the parameters of g_w , at both meta-training and evaluation is modified as follows:

$$w_t = \begin{cases} w_{t-1}^i - \frac{1}{\text{FAN_IN}} \cdot \left(\alpha_w \lambda_1 d_\phi^i \exp \left(\lambda_2 m_\phi^i \right) \right) & \text{if } w^i \text{ is part of a hidden layer} \\ w_{t-1}^i - \alpha_w \lambda_1 d_\phi^i \exp \left(\lambda_2 m_\phi^i \right) & \text{otherwise.} \end{cases} \quad (1)$$

We now show that this can lead to a maximal update Parametrization, following the analysis of [36, Appendix J.2.1] which studies the initial optimization step. For our analysis, we consider a simplified input set for f_ϕ which takes as input only the gradient while producing an update for each layer. Note that this analysis extends naturally to other first-order quantities.

Proposition 1. *Assume that the LO f_ϕ is continuous around 0. Then, if $f_\phi(0) \neq 0$, the update, initialization, and pre-activation multipliers above are necessary to obtain a Maximal Update Parametrization.*

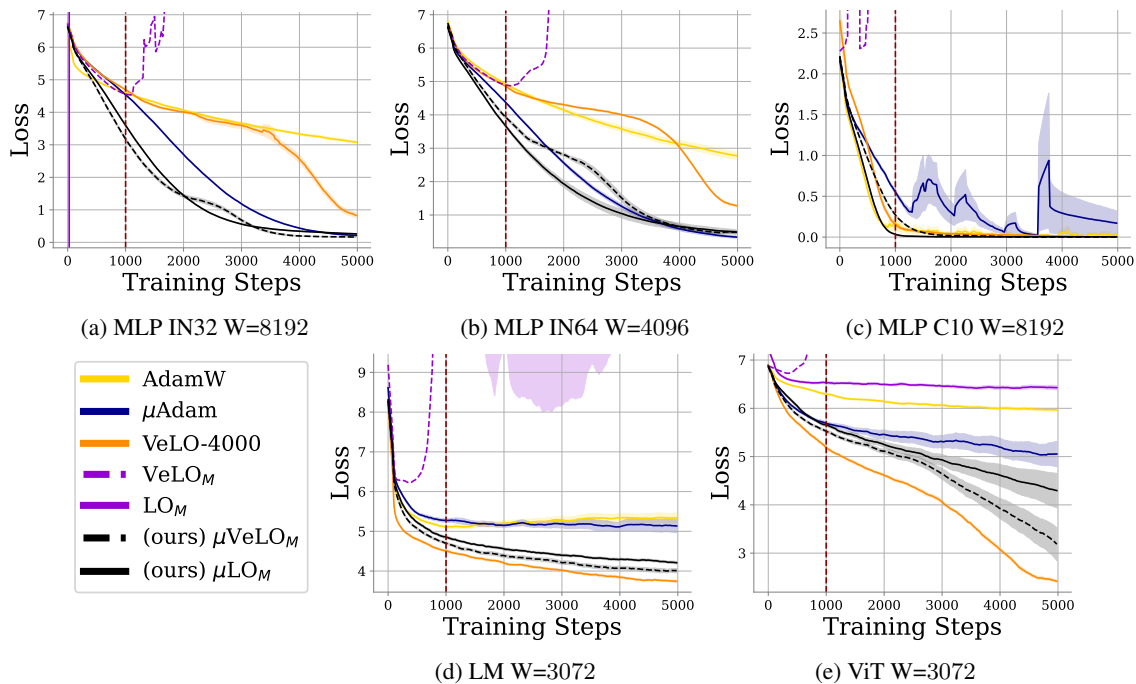


Figure 2: **Evaluating generalization to wider networks for different tasks.** Tasks Our optimizers are meta-trained for 1000 inner steps (dotted red line), therefore, any optimization beyond 1000 steps is considered out-of-distribution. We plot average training loss over 5 seeds with standard error bars. We observe that μ LO_M and μ VeLO_M generalize smoothly to longer unrolls and all unseen tasks, unlike their SP counterparts which diverge or fail to make progress. μ LOs even surpass or match the performance of VeLO in subfigures (a), (b), and (c)). Moreover, they also substantially best the well-tuned hand-designed baselines on LM and ViT tasks (subfigures (d) and (e)) and best or match the best performing hand-designed optimizer in subfigures (a),(b), and (c).

3. Empirical evaluation

We use a suite of optimization tasks to evaluate meta-generalization properties of our μ LOs vs tuned μ Adam [36] and baseline LOs, including pretrained VeLO [20] (the strongest publicly available pre-trained learned optimizer). We focus on evaluating generalization to wider networks, however, we also verify if μ LOs can improve generalization to longer training horizons and deeper networks. Extended details of our experimental setup and baselines are provided in the appendix section E.

3.1. Meta-generalization to wider networks

In this section, we briefly outline the results showing the improved generalization of μ LOs to wider networks. Firstly, Figure 1 illustrates the improvements of μ LOs over standard parameterization learned optimizers: the latter do not generalize to networks wider than those seen at meta-training time. Furthermore, figure 2 reports the training curves of different optimizers on the largest width tasks in our suite. Despite training for $5\times$ longer than the maximum meta-training unroll length, our μ LOs are capable of smoothly decreasing the loss for the largest out-of-distribution tasks in our suite. In contrast, the strong SP LO baselines diverge by 1000 steps (subfigures (a),(b),(c),(d)), or fail to decrease the training loss (subfigure (e)). Our μ LOs also substantially best the well-tuned

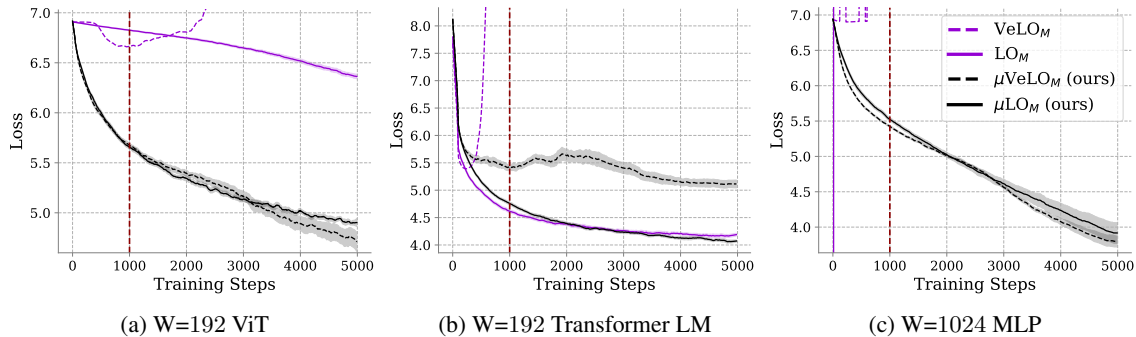


Figure 3: **Evaluating generalization capabilities of μ LOs to deeper networks.** The figures report the performance of learned optimizers for training depth-16 ViTs for image classification, Transformers for language modeling, and MLPs for image classification. We plot average training loss over 5 seeds with standard error bars. In each case, μ LOs show improved generalization and performance when compared to their SP counterparts.

hand-designed baselines on LM and ViT tasks (subfigures (d) and (e)) and best or match the best performing hand-designed optimizer in subfigures (a),(b), and (c). When comparing with VeLO-4000, we observe that our μ LOs substantially outperform VeLO in subfigures (a),(b), and μ LO_M outperforms VeLO-4000 in subfigure (c). Notably, VeLO-4000 only outperforms the μ LOs on tasks they did not see during meta-training.

3.2. Meta-generalization to deeper networks

In this section, we evaluate LO meta-generalization to deeper networks. Specifically, we increase the number of layers used in MLP, ViT, and LM tasks from 3 to 16, while being sure to select models that have widths within the meta-training range (128 – 1024) to avoid confounding the results. Figure 3 reports the performance of our multi-task learned optimizers on deeper networks. We observe that both μ LO_M and μ VeLO_M optimize stably throughout and generally outperform their counterparts, LO_M and VeLO_M, by the end of training on each task, despite being meta-trained on MLPs of exactly the same depth. Moreover, LO_M immediately diverges when optimizing the deep MLP while μ LO_M experience no instability. Similarly, VeLO_M diverges on ViTs and Transformers, while μ VeLO_M performs well, especially on ViTs. This is remarkable as, unlike width, there is no theoretical justification for μ P’s benefit to deeper networks. We hypothesize that μ P’s stabilizing effect on the optimizee’s activations leads to this improvement generalization.

3.3. Meta-generalization to longer training horizons

In this subsection, we empirically evaluate the capability of μ LOs to generalize to much longer training horizons than those seen during meta-training. Specifically, we use μ LO_M and LO_M as well as μ VeLO_M and VeLO_M to train three networks with width $w = 1024$: a 3-layer MLP, ViT on $32 \times 32 \times 3$ ImageNet and a 3-layer Transformer for autoregressive language modeling on LM1B. Each model is trained for 25,000 steps ($25 \times$ the longest unroll seen at meta-training time). Figure 4 reports the training loss averaged over 5 random seeds. We observe that μ LO_M and μ VeLO_M stably decrease training loss over time for each task, while LO_M and VeLO_M fail to decrease training loss (a), decreases it but suffers instabilities (b), or diverges after 8000 steps (c). These results suggest that generalization to longer training horizons is another benefit of using μ P with learned optimizers.

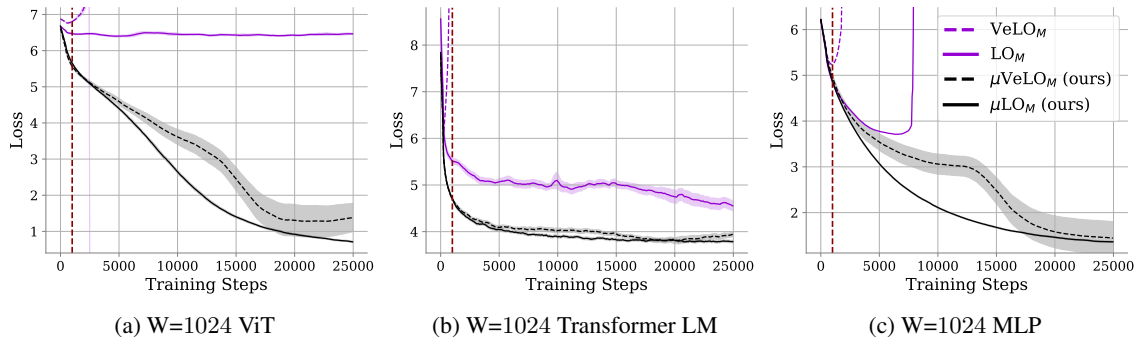


Figure 4: **Evaluating generalization capabilities of μ LOs to longer training horizons.** We plot average training loss over 5 seeds with standard error bars. All optimizers are meta-trained for 1000 steps of training (dotted red line), therefore, any optimization beyond 1000 steps is considered out-of-distribution. We observe that μ LOs seamlessly generalize to training horizons $25\times$ longer than meta-training. In contrast, the best performing SP LO fails to decrease training loss (a), decreases it but suffers instabilities (b), or diverges after 8000 steps (c).

4. Limitations

While we have conducted a systematic empirical study and shown strong results within the scope of our study, there are some of limitations of our work. Specifically, (1) we do not meta-train on tasks other than MLPs for image classification and we do not provide evaluation of models wider than 8192 for MLPs and 3072/12288 (hidden/FFN size) for transformers due to computational constraints in our academic environment.

5. Conclusion

We have demonstrated that applying our μ LO meta-training recipe produces optimizers with substantially improved meta-generalization properties when compared to strong baselines from previous work. Remarkably, our μ LOs even surpass VeLO-4000 (meta-trained for 4000 TPU months) on wider versions of in-distribution tasks. Moreover, our experiments also show that μ LOs meta-trained with our recipe generalize better to wider and deeper out-of-distribution tasks than their SP counterparts. When evaluated on much longer training tasks, we observe that μ LOs have a stabilizing effect, enabling meta-generalization to much longer unrolls ($25\times$ maximum meta-training unroll length). All of the aforementioned benefits of μ LOs come at *zero* extra computational cost compared to SP LOs. Our results outline a promising path forward for low-cost meta-training of learned optimizers that can generalize to large unseen tasks.

Acknowledgements

We acknowledge support from Mila-Samsung Research Grant, FRQNT New Scholar [E.B.], the FRQNT Doctoral (B2X) scholarship [B.T.], the Canada CIFAR AI Chair Program [I.R.], the French Project ANR-21-CE23-0030 ADONIS [E.O.], and the Canada Excellence Research Chairs Program in Autonomous AI [I.R.]. We also acknowledge resources provided by Compute Canada, Calcul Québec, and Mila.

References

- [1] Diogo Almeida, Clemens Winter, Jie Tang, and Wojciech Zaremba. A generalizable approach to learning optimizers. *arXiv preprint arXiv:2106.00958*, 2021.
- [2] Brandon Amos. Tutorial on amortized optimization for learning to optimize over continuous domains. *arXiv e-prints*, pages arXiv–2202, 2022.
- [3] Marcin Andrychowicz, Misha Denil, Sergio Gomez, Matthew W Hoffman, David Pfau, Tom Schaul, Brendan Shillingford, and Nando De Freitas. Learning to learn by gradient descent by gradient descent. *Advances in neural information processing systems*, 29, 2016.
- [4] Jacob Buckman, Danijar Hafner, George Tucker, Eugene Brevdo, and Honglak Lee. Sample-efficient reinforcement learning with stochastic ensemble value expansion. In Samy Bengio, Hanna M. Wallach, Hugo Larochelle, Kristen Grauman, Nicolò Cesa-Bianchi, and Roman Garnett, editors, *Advances in Neural Information Processing Systems 31: Annual Conference on Neural Information Processing Systems 2018, NeurIPS 2018, December 3-8, 2018, Montréal, Canada*, pages 8234–8244, 2018. URL <https://proceedings.neurips.cc/paper/2018/hash/f02208a057804ee16ac72ff4d3cec53b-Abstract.html>.
- [5] Ciprian Chelba, Tomas Mikolov, Mike Schuster, Qi Ge, Thorsten Brants, and Phillipp Koehn. One billion word benchmark for measuring progress in statistical language modeling. *CoRR*, abs/1312.3005, 2013. URL <http://arxiv.org/abs/1312.3005>.
- [6] Tianlong Chen, Weiyi Zhang, Zhou Jingyang, Shiyu Chang, Sijia Liu, Lisa Amini, and Zhangyang Wang. Training stronger baselines for learning to optimize. *Advances in Neural Information Processing Systems*, 33:7332–7343, 2020.
- [7] Tianlong Chen, Xiaohan Chen, Wuyang Chen, Zhangyang Wang, Howard Heaton, Jialin Liu, and Wotao Yin. Learning to optimize: A primer and a benchmark. *The Journal of Machine Learning Research*, 23(1):8562–8620, 2022.
- [8] N. Dalal and B. Triggs. Histograms of oriented gradients for human detection. In *2005 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR'05)*, volume 1, pages 886–893 vol. 1, 2005. doi: 10.1109/CVPR.2005.177.
- [9] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*, 2020.
- [10] Katie E. Everett, Lechao Xiao, Mitchell Wortsman, Alexander A. Alemi, Roman Novak, Peter J. Liu, Izzeddin Gur, Jascha Sohl-Dickstein, Leslie Pack Kaelbling, Jaehoon Lee, and Jeffrey Pennington. Scaling exponents across parameterizations and optimizers. In *Forty-first International Conference on Machine Learning, ICML 2024, Vienna, Austria, July 21-27, 2024*. OpenReview.net, 2024. URL <https://openreview.net/forum?id=0ksNeD1SJT>.

- [11] James Harrison, Luke Metz, and Jascha Sohl-Dickstein. A closer look at learned optimization: Stability, robustness, and inductive biases. *Advances in Neural Information Processing Systems*, 35:3758–3773, 2022.
- [12] Charles-Étienne Joseph, Benjamin Thérien, Abhinav Moudgil, Boris Knyazev, and Eugene Belilovsky. Can we learn communication-efficient optimizers? *CoRR*, abs/2312.02204, 2023. URL <https://doi.org/10.48550/arXiv.2312.02204>.
- [13] Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization, 2017.
- [14] Taku Kudo and John Richardson. Sentencepiece: A simple and language independent subword tokenizer and detokenizer for neural text processing. In Eduardo Blanco and Wei Lu, editors, *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing, EMNLP 2018: System Demonstrations, Brussels, Belgium, October 31 - November 4, 2018*, pages 66–71. Association for Computational Linguistics, 2018. URL <https://doi.org/10.18653/v1/d18-2012>.
- [15] Haochuan Li, Alexander Rakhlin, and Ali Jadbabaie. Convergence of adam under relaxed assumptions. In Alice Oh, Tristan Naumann, Amir Globerson, Kate Saenko, Moritz Hardt, and Sergey Levine, editors, *Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023*, 2023. URL http://papers.nips.cc/paper_files/paper/2023/hash/a3cc50126338b175e56bb3cad134db0b-Abstract-Conference.html.
- [16] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. In *7th International Conference on Learning Representations, ICLR 2019, New Orleans, LA, USA, May 6-9, 2019*. OpenReview.net, 2019. URL <https://openreview.net/forum?id=Bkg6RiCqY7>.
- [17] David G. Lowe. Distinctive image features from scale-invariant keypoints. *Int. J. Comput. Vis.*, 60(2):91–110, 2004. URL <https://doi.org/10.1023/B:VISI.0000029664.99615.94>.
- [18] Luke Metz, Niru Maheswaranathan, Jeremy Nixon, Daniel Freeman, and Jascha Sohl-Dickstein. Understanding and correcting pathologies in the training of learned optimizers. In *International Conference on Machine Learning*, pages 4556–4565. PMLR, 2019.
- [19] Luke Metz, C. Daniel Freeman, James Harrison, Niru Maheswaranathan, and Jascha Sohl-Dickstein. Practical tradeoffs between memory, compute, and performance in learned optimizers, 2022.
- [20] Luke Metz, James Harrison, C Daniel Freeman, Amil Merchant, Lucas Beyer, James Bradbury, Naman Agrawal, Ben Poole, Igor Mordatch, Adam Roberts, et al. Velo: Training versatile learned optimizers by scaling up. *arXiv preprint arXiv:2211.09760*, 2022.
- [21] Yurii E. Nesterov and Vladimir G. Spokoiny. Random gradient-free minimization of convex functions. *Found. Comput. Math.*, 17(2):527–566, 2017. URL <https://doi.org/10.1007/s10208-015-9296-2>.

- [22] Paavo Parmas, Carl Edward Rasmussen, Jan Peters, and Kenji Doya. PIPPS: flexible model-based policy search robust to the curse of chaos. In Jennifer G. Dy and Andreas Krause, editors, *Proceedings of the 35th International Conference on Machine Learning, ICML 2018, Stockholmsmässan, Stockholm, Sweden, July 10-15, 2018*, volume 80 of *Proceedings of Machine Learning Research*, pages 4062–4071. PMLR, 2018. URL <http://proceedings.mlr.press/v80/parmas18a.html>.
- [23] Isabeau Premont-Schwarz, Jaroslav Vitkuu, and Jan Feyereisl. A simple guard for learned optimizers. *arXiv preprint arXiv:2201.12426*, 2022.
- [24] Herbert E. Robbins. A stochastic approximation method. *Annals of Mathematical Statistics*, 22:400–407, 1951. URL <https://api.semanticscholar.org/CorpusID:16945044>.
- [25] Jürgen Schmidhuber. Learning to control fast-weight memories: An alternative to dynamic recurrent networks. *Neural Computation*, 4(1):131–139, 1992.
- [26] James C. Spall. Adaptive stochastic approximation by the simultaneous perturbation method. *IEEE Trans. Autom. Control.*, 45(10):1839–1853, 2000. URL <https://doi.org/10.1109/TAC.2000.880982>.
- [27] Sebastian Thrun and Lorien Pratt. *Learning to learn*. Springer Science & Business Media, 2012.
- [28] Paul Vicol. Low-variance gradient estimation in unrolled computation graphs with es-single. In *International Conference on Machine Learning*, pages 35084–35119. PMLR, 2023.
- [29] Paul Vicol, Luke Metz, and Jascha Sohl-Dickstein. Unbiased gradient estimation in unrolled computation graphs with persistent evolution strategies. In Marina Meila and Tong Zhang, editors, *Proceedings of the 38th International Conference on Machine Learning, ICML 2021, 18-24 July 2021, Virtual Event*, volume 139 of *Proceedings of Machine Learning Research*, pages 10553–10563. PMLR, 2021. URL <http://proceedings.mlr.press/v139/vicol21a.html>.
- [30] Olga Wichrowska, Niru Maheswaranathan, Matthew W Hoffman, Sergio Gomez Colmenarejo, Misha Denil, Nando Freitas, and Jascha Sohl-Dickstein. Learned optimizers that scale and generalize. In *International conference on machine learning*, pages 3751–3760. PMLR, 2017.
- [31] Greg Yang. Tensor programs I: wide feedforward or recurrent neural networks of any architecture are gaussian processes. *CoRR*, abs/1910.12478, 2019. URL <http://arxiv.org/abs/1910.12478>.
- [32] Greg Yang. Tensor programs II: neural tangent kernel for any architecture. *CoRR*, abs/2006.14548, 2020. URL <https://arxiv.org/abs/2006.14548>.
- [33] Greg Yang. Tensor programs III: neural matrix laws. *CoRR*, abs/2009.10685, 2020. URL <https://arxiv.org/abs/2009.10685>.

- [34] Greg Yang and Edward J. Hu. Tensor programs IV: feature learning in infinite-width neural networks. In Marina Meila and Tong Zhang, editors, *Proceedings of the 38th International Conference on Machine Learning, ICML 2021, 18-24 July 2021, Virtual Event*, volume 139 of *Proceedings of Machine Learning Research*, pages 11727–11737. PMLR, 2021. URL <http://proceedings.mlr.press/v139/yang21c.html>.
- [35] Greg Yang and Etai Littwin. Tensor programs ivb: Adaptive optimization in the infinite-width limit. *CoRR*, abs/2308.01814, 2023. URL <https://doi.org/10.48550/arXiv.2308.01814>.
- [36] Greg Yang, Edward J. Hu, Igor Babuschkin, Szymon Sidor, David Farhi, Jakub Pachocki, Xiaodong Liu, Weizhu Chen, and Jianfeng Gao. Tensor programs v: Tuning large neural networks via zero-shot hyperparameter transfer. In *NeurIPS 2021*, March 2022. URL <https://www.microsoft.com/en-us/research/publication/tuning-large-neural-networks-via-zero-shot-hyperparameter-transfer/>.
- [37] Greg Yang, Dingli Yu, Chen Zhu, and Soufiane Hayou. Tensor programs VI: feature learning in infinite depth neural networks. In *The Twelfth International Conference on Learning Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024*. OpenReview.net, 2024. URL <https://openreview.net/forum?id=17pVDnpwwl>.
- [38] Junjie Yang, Tianlong Chen, Mingkang Zhu, Fengxiang He, Dacheng Tao, Yingbin Liang, and Zhangyang Wang. Learning to generalize provably in learning to optimize. In *International Conference on Artificial Intelligence and Statistics*, pages 9807–9825. PMLR, 2023.

Appendix A. Table of Contents**Contents**

1	Introduction	1
2	μ-parametrization for Learned Optimizers	3
3	Empirical evaluation	4
3.1	Meta-generalization to wider networks	4
3.2	Meta-generalization to deeper networks	5
3.3	Meta-generalization to longer training horizons	5
4	Limitations	6
5	Conclusion	6
A	Table of Contents	11
B	Proof of proposition 1	12
C	Related Work	12
C.1	Learned optimization.	12
C.2	Generalization in LOs.	12
C.3	Maximal Update Parametrization.	13
D	Learned Optimization Background	13
D.1	μ LO Meta-training Recipe	14
E	Experimental Setup	15
E.1	Optimizers in our study	15
E.2	List of Evaluation Tasks	16
F	Extended Results	18
F.1	Evaluating pre-activation stability	18
G	Hand Designed Optimizer Hyperparameter Tuning	18
G.1	Tuning μ Adam	18
G.2	Tuning AdamW	18
H	Meta-training with μLOs	19
I	Features of the learned optimizer	21
J	Coordinate evolution of MLP layers in μP for different optimizers	22

Appendix B. Proof of proposition 1

Proposition 1. *Assume that the LO f_ϕ is continuous around 0. Then, if $f_\phi(0) \neq 0$, the update, initialization, and pre-activation multipliers above are necessary to obtain a Maximal Update Parametrization.*

Proof The update produced by f_ϕ is denoted ΔW and we write ∇W the corresponding gradient, so that $\Delta W = f_\phi(\nabla W)$. For the sake of simplicity, n will be the output size and d the feature input size of our neural network. Our goal is to satisfy the desiderata of [36, Appendix J.2]. We assume our initialization follows Initialization- μ in Sec 3. Overall, our goal is to study strategy so that if $x_i = \Theta(1)$, then one needs to renormalize/initialize so that $(Wx)_i = \Theta(1)$ while $((W + \Delta W)x)_i = \Theta(1)$ so that the update is as large as possible. Note that given the assumptions on f , if $x = \Theta(\frac{1}{n})$, then $f(x) = \Theta(1)$.

Output weights. Here, if input x has some $\Theta(1)$ coordinates, we initialize $W = (w_i)_{i \leq n}$ with weights of variance 1 (which is necessary) and rescale the preactivations with $\frac{1}{n}$. For the update, we thus have that ∇W scales (coordinate wise) as $\Theta(\frac{1}{n})$ and we do not rescale the LR, given that $f_\phi(\nabla W)$ will also have coordinates in $\Theta(1)$.

Hidden weights. Now, for the update, we observe that the gradient ∇W has some coordinates which scale as $\Theta(\frac{1}{n})$, due to the output renormalization choice. Thus, the LO $f_\phi(\nabla W)$ satisfies that $f(\nabla W) = \Theta(1)$, given that f_ϕ is continuous in 0 and satisfies $f_\phi(0) \neq 0$. Thus for the update, we need to use $\Delta W = \frac{1}{n} f_\phi(\nabla W)$ so that $\Delta W x$ is coordinate wise bounded.

Input weights. In this case, the gradient has coordinates which already scale in $\Theta(\frac{1}{n})$ (due to the output renormalization) and there is no need to rescale the LR.

Appendix C. Related Work

C.1. Learned optimization.

While research on learned optimizers (LOs) spans several decades [2, 7, 25, 27], our work is primarily related to the recent meta-learning approaches utilizing efficient per-parameter optimizer architectures of [19]. Unlike prior work [3, 6, 30], which computes meta-gradients (the gradients of the learned optimizer) using backpropagation, [19] use Persistent Evolutionary Strategies (PES) [29], a truncated variant of evolutionary strategies (ES) [4, 21, 22]. ES improves meta-training of LOs by having more stable meta-gradient estimates compared to backpropagation through time, especially for longer sequences (i.e. long parameter update unrolls inherent in meta-training) [18]. PES and most recently ES-Single [28] are more efficient and accurate variants of ES, among which PES is more well-established in practice making it a favourable approach to meta-training.

C.2. Generalization in LOs.

One of the critical issues in LOs is generalization in the three main aspects [2, 7]: (1) optimize novel tasks (often referred to as *meta-generalization*); (2) optimize for more iterations than the maximum unroll length used in meta-training; (3) avoid overfitting on the training set. Among these, (3) has been extensively addressed using different approaches, such as meta-training on the validation set objective [18], adding extra-regularization terms [11], parameterizing LOs as hyperparameter controllers [1] and introducing flatness-aware regularizations [38]. The regularization terms [11, 38] often alleviate issue (2) as a byproduct. However, meta-generalization (1) has remained a more

difficult problem. One approach to tackle this problem is to meta-train LOs on thousands of tasks [20]. However, this approach is extremely expensive and does not address the issue in a principled way leading to poor meta-generalization in some cases, e.g. when the optimization task includes much larger networks. Alternatively, [23] introduced Loss-Guarded L2O (LGL2O) that switches to Adam/SGD if the LO starts to diverge improving meta-generalization. However, this approach needs tuning Adam/SGD and requires additional computation (e.g. for loss check) limiting (or completely diminishing in some cases) the benefits of the LO. In this work, we study aspects (1) and (2) of LO generalization, demonstrating how existing SP LOs generalize poorly across these dimensions and showing how one can apply μ P to learned optimizers to substantially improve generalization in both these aspects.

C.3. Maximal Update Parametrization.

First proposed by [34], the Maximal Update Parametrization is the unique stable abc-Parametrization where every layer learns features. The parametrization was derived for adaptive optimizers by [35] and was applied by [36] to enable zero-shot hyperparameter transfer, constituting the first practical application of the tensor programs series of papers. Earlier works in the *tensor programs series* build the mathematical foundation that led to the discovery of μ P. [31] shows that many modern neural networks with randomly initialized weights and biases are Gaussian Processes, providing a language, called Netsor, to formalize neural network computations. [32] focuses on neural tangent kernels (NTK), proving that as a randomly initialized network’s width tends to infinity, its NTK converges to a deterministic limit. [33] shows that randomly initialized network’s pre-activations become independent of its weights when its width tends to infinity. Most recently, in tensor programs VI, [37] propose Depth- μ P, a parameterization allowing for hyperparameter transfer in infinitely deep networks. However, Depth- μ P is only valid for residual networks with a block depth of 1, making it unusable for most practical architectures (e.g., transformers, resnets, etc.). For these reasons, we do not study Depth- μ P herein. Building on the latest works studying width μ P [35, 36], in this work, we show that μ P can be extended to the case of learned optimizers and empirically evaluate its benefits in this setting.

Appendix D. Learned Optimization Background

A standard approach to learning optimizers [19] is to solve the following meta-learning problem:

$$\min_{\phi} \mathbb{E}_{(\mathcal{D}, \mathbf{w}_0) \sim \mathcal{T}} \left[\mathbb{E}_{(X, Y) \sim \mathcal{D}} \left[\frac{1}{T} \sum_{t=0}^{T-1} \mathcal{L}(X, Y; f_{\phi}(\mathbf{u}_t), \mathbf{w}_t) \right] \right], \quad (2)$$

where \mathcal{T} is a distribution over optimization tasks defined as pairs of dataset \mathcal{D} and initial weights \mathbf{w}_0 associated with a particular neural architecture (we refer to this network as the *optimizee*), ϕ represents the weights of the learned optimizer, f_{ϕ} , that takes gradient-based features \mathbf{u}_t as input. Finally, \mathcal{L} is the loss function used to train the optimizee. T is the length of the unroll which we write as a fixed quantity for simplicity. In our experiments, during meta-optimization, T is varied according to a truncation schedule [19]. A clear goal of the learned optimization community is not only learning to solve optimization problems over \mathcal{T} , but also to apply the learned optimizer, f_{ϕ} , more generally to unobserved datasets and architectures. This *transfer* to new tasks is referred to as meta-generalization. This problem can be seen as a generalization of the zero-shot hyperparameter transfer problem considered in [36]; for instance, when the optimizer is a hand-designed method such as SGD or Adam and ϕ represents optimization hyper-parameters such as the learning rate.

Gradient descent is a standard approach to solving equation 2. However, estimating the meta-gradients via backpropagation for very long unrolls is known to be noisy [18]. Instead, gradients are estimated using evolution strategies [4, 21, 22]. Evolution strategies work by sampling perturbations to the LO’s weights (similar to SPSA [26]), unrolling an optimization trajectory for each perturbation, and estimating gradients with respect to evaluations of the meta-objective (usually the loss of the network being optimized, see eq. 2). In contrast to ES, which estimates one gradient per full unroll, PES [29] allows estimating unbiased gradients at many points (called truncations) during the full unroll. This allows updating the optimizer’s parameters more often during meta-training. We use PES to estimate meta-gradients in our experiments.

Learned optimizer features u_t are constructed based on momentum, second-order momentum, and adafactor values as in [19], with the full list of features described in the (Table 6 of the Appendix). In our experiments, the architectures of our f_ϕ are similar to **small_fc_lopt** of [19] and **VeLO** of [20] except that their dimensions differ slightly (see sec. H for details). f_ϕ has three outputs m , d , and α , the magnitude, scale, and learning rate of the update respectively. For **small_fc_lopt**, $\alpha_w = 1$ always, α_w is produced by the tensor-level LSTM for VeLO. The standard LO update is given as

$$w_t = w_{t-1} - \alpha_w \lambda_1 d_\phi \exp(\lambda_2 m_\phi), \tag{3}$$

where λ_1 and λ_2 are constant values of 0.001 to bias initial step sizes towards being small.

D.1. μ LO Meta-training Recipe

μ P for hand-designed optimizers involves tuning the optimizer on a small width version of the target architecture and transferring the hyperparameters to the larger width target model [36]. While μ -transfer makes hyperparameter search for large models tractable, it has the following limitations: (1) the smaller scale hyperparameter search suffers from increased complexity as it requires sweeping various multipliers in addition to the standard hyperparameters, (2) tuning the hyperparameters on too small of a model may result in sub-optimal hyperparameters for the largest models, (3) [36] recommend repeating the procedure for every new task/dataset. Meta-training flexible μ -parametrized learned optimizers can address these limitations. Due to their flexible functional forms (as opposed to just a learning rate hyperparameter), μ LOs can learn to optimize networks in μ P without tuning multipliers (we set all multipliers to 1 in our experiments). Therefore, by training our μ LOs with fixed multipliers on *multiple tasks* that are large enough to admit strong transfer but still tractable and reusing them on new tasks, we address (1), (2), and (3) by amortizing the tuning cost during the optimizer meta-training stage. However, it should be noted that while the μ LO

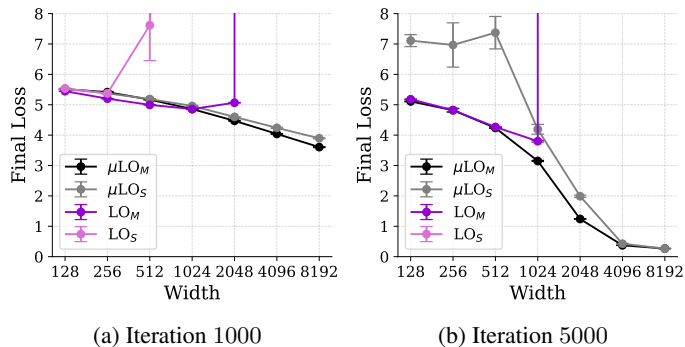


Figure 5: μ LO_S underperforms μ LO_M as width and training steps increase. Each point is the average training loss over 5 seeds at iterations 1000 (a) or 5000 (b). Error bars report standard error.

Table 1: **Meta-training configurations of LOs and baselines in our empirical evaluation.**

Identifier	Type	Architecture	Optimizee Par.	Meta-Training / Tuning Task(s)
μ LO _S	Ours	small_fc_lopt [19]	μ LO Sec. 2	ImageNet classification, 3-Layer MLP, width $\in \{128\}$
μ LO _M	Ours	small_fc_lopt [19]	μ LO Sec. 2	ImageNet classification, 3-Layer MLP, width $\in \{128, 512, 1024\}$
μ VeLO _M	Ours	VeLO [20]	μ LO Sec. 2	ImageNet classification, 3-Layer MLP, width $\in \{128, 512, 1024\}$
LO _S	LO Baseline	small_fc_lopt [19]	SP	ImageNet classification, 3-Layer MLP, width $\in \{128\}$
LO _M	LO Baseline	small_fc_lopt [19]	SP	ImageNet classification, 3-Layer MLP, width $\in \{128, 512, 1024\}$
VeLO _M	LO Baseline	VeLO [20]	SP	ImageNet classification, 3-Layer MLP, width $\in \{128, 512, 1024\}$
VeLO-4000	Oracle LO Baseline	VeLO [20]	SP	We refer the reader to [20, Appendix C.2]
μ Adam	Baseline	–	μ P Adam	ImageNet classification, 3-Layer MLP, width $\in \{1024\}$
AdamW	Baseline	–	SP	ImageNet classification, 3-Layer MLP, width $\in \{1024\}$

framework allows for meta-generalization to unseen new tasks (unlike μ -transfer), a μ LO that relies on meta-generalization for transfer to new tasks should expect to be outperformed by the μ LO that also meta-trains on a small version of that task.

To verify the effectiveness of this multi-task strategy for learned optimizers, we compare μ LO_S, trained on a single small task (see Tab. 1), to μ LO_M, trained on 3 small tasks of the different width (see Tab. 1), in figure 5. When training for 1000 steps (meta-training unroll length), we observe that μ LO_M outperforms μ LO_S as the width of the model is increased (Fig. 5 (a)). Moreover, we observe that there is a discrepancy in performance between both models after 5000 steps (Fig. 5 (b)), showing that meta-training with multiple tasks of different widths has benefits for generalization to longer unrolls in addition to improved generalization to larger optimizees. Given the improved generalization of μ LO_M compared to μ LO_S, we adopt the multiple-width single-task meta-training recipe as part of our method. Subsequent experiments (e.g., figures 1 and 2) will show that it is also effective for meta-training μ VeLO.

Appendix E. Experimental Setup

We use a suite of optimization tasks of varying width to evaluate meta-generalization properties of our μ LOs vs tuned μ Adam [36], SP AdamW [16], and baseline SP LOs. We also include pre-trained VeLO [20] as an oracle which we denote as VeLO-4000. Meta-trained for 4000 TPUv4 months, it is the strongest publicly available pre-trained learned optimizer. We focus on evaluating generalization to wider networks, however, we also establish the generalization properties of μ LOs to longer training horizons and deeper networks. Please note that while μ LOs inherit the theoretical properties of μ P for width scaling, our findings with respect to longer training and deeper networks are purely empirical.

E.1. Optimizers in our study

Baseline LOs and μ LOs. The meta-training configuration of each learned optimizer is summarized in Table 1. Each learned optimizer (ours and the baselines) in our empirical evaluation is meta-trained using the multiple-width single-task meta-training recipe proposed in section D.1. The baseline sheds light on whether simply varying the SP optimizee width during meta-training is enough to achieve generalization of the LO to wider networks in SP. During meta-training, we set the inner problem length to be 1000 iterations. Therefore, any optimization beyond this length is considered out-of-distribution. For all meta-training and hyperparameter tuning details, including ablation experiments, see section H of the appendix.

μ Adam μ Adam is a strong hand-designed μ P baseline. It follows the [36] Adam μ -parametrization and does not use weight decay as this is incompatible with μ P. It is tuned on the largest meta-training

task seen by our learned optimizers (Table 1). We tune the learning rate and three multipliers: input multiplier, output multiplier, and the hidden learning rate multiplier. These multipliers correspond to adding a tunable constant to the pre-activation multiplier for input weights, the pre-activation multiplier for output weights, and the Adam LR for hidden weights. More details about the grid search over 500 configurations are provided in Section G.1 of the appendix.

AdamW AdamW [16] is a strong hand-designed SP baseline. It is tuned on the largest meta-training task seen by our learned optimizers (Table 1). We tune the learning rate, β_1, β_2 , and the weight decay. More details about the grid search over 500 configurations are provided in Section G.1 of the appendix.

Pre-trained VeLO (VeLO-4000). VeLO [20] is a learned optimizer that was meta-trained on a curriculum of progressively more expensive meta-training tasks for a total of 4000 TPU months. These tasks include but are not limited to image classification with MLPs, ViTs, ConvNets, and ResNets; compression with MLP auto-encoders; generative modeling with VAEs; and language modeling with transformers and recurrent neural networks. During meta-training, VeLO-4000 unrolls inner problems for up to 20k steps ($20\times$ ours); the final model was then fine-tuned on tasks with up to 200k steps of optimization. VeLO-4000, therefore, represents the strongest baseline in our empirical evaluation and we consider it to be an oracle.

Is VeLO-4000 a fair baseline? While we believe the comparison is important given the relevance of our results to scaling up LOs, we highlight that the comparison will unfairly advantage VeLO-4000 as all tasks in our suite fall within its meta-training distribution and VeLO-4000 was meta-trained on inner unroll horizons well beyond those we evaluate. Thus, when comparing our LOs to VeLO-4000, it is important to keep in mind that ours are meta-trained with only 0.004% of VeLO-4000’s compute budget.

Evaluation tasks. Our evaluation suite includes 35 tasks spanning image classification (CIFAR-10, ImageNet) using MLPs and Vision Transformers (ViTs) [9] and autoregressive language modeling with a decoder-only transformer on LM1B [5]. To create the tasks, we further vary image size (for image classification), width, and depth of the optimizee network, and the number of optimization steps. See Table 2 of the appendix for an extended description of all the tasks.

E.2. List of Evaluation Tasks

Table 2 reports the configuration of different testing tasks used to evaluate our optimizers. We note that we do not augment the ImageNet datasets we use in any way except for normalizing the images. We tokenize LM1B using a sentence piece tokenizer[14] with 32k vocabulary size. All evaluation tasks are run on A6000 or A100 GPUs for 5 random seeds. Each individual run takes less than 2 hours to complete on a single 80GB A100.

Table 2: **Meta-testing settings.** We report the optimization tasks we will use to evaluate the LOs of Table 1.

Identifier	Dataset	Model	Depth	Width	Attn. Heads	FFN Size	Batch Size	Sequence Length
$\text{IN32T}_{(3,128)}^{\text{MLP}}$	$32 \times 32 \times 3$ ImageNet	MLP	3	128	–	–	4096	–
$\text{IN32T}_{(3,256)}^{\text{MLP}}$	$32 \times 32 \times 3$ ImageNet	MLP	3	256	–	–	4096	–
$\text{IN32T}_{(3,512)}^{\text{MLP}}$	$32 \times 32 \times 3$ ImageNet	MLP	3	512	–	–	4096	–
$\text{IN32T}_{(3,1024)}^{\text{MLP}}$	$32 \times 32 \times 3$ ImageNet	MLP	3	1024	–	–	4096	–
$\text{IN32T}_{(3,2048)}^{\text{MLP}}$	$32 \times 32 \times 3$ ImageNet	MLP	3	2048	–	–	4096	–
$\text{IN32T}_{(3,4096)}^{\text{MLP}}$	$32 \times 32 \times 3$ ImageNet	MLP	3	4096	–	–	4096	–
$\text{IN32T}_{(3,8192)}^{\text{MLP}}$	$32 \times 32 \times 3$ ImageNet	MLP	3	8192	–	–	4096	–
$\text{IN64T}_{(3,128)}^{\text{MLP}}$	$64 \times 64 \times 3$ ImageNet	MLP	3	128	–	–	4096	–
$\text{IN64T}_{(3,256)}^{\text{MLP}}$	$64 \times 64 \times 3$ ImageNet	MLP	3	256	–	–	4096	–
$\text{IN64T}_{(3,512)}^{\text{MLP}}$	$64 \times 64 \times 3$ ImageNet	MLP	3	512	–	–	4096	–
$\text{IN64T}_{(3,1024)}^{\text{MLP}}$	$64 \times 64 \times 3$ ImageNet	MLP	3	1024	–	–	4096	–
$\text{IN64T}_{(3,2048)}^{\text{MLP}}$	$64 \times 64 \times 3$ ImageNet	MLP	3	2048	–	–	4096	–
$\text{IN64T}_{(3,4096)}^{\text{MLP}}$	$64 \times 64 \times 3$ ImageNet	MLP	3	4096	–	–	4096	–
$\text{C10T}_{(3,128)}^{\text{MLP}}$	$32 \times 32 \times 3$ Cifar-10	MLP	3	128	–	–	4096	–
$\text{C10T}_{(3,256)}^{\text{MLP}}$	$32 \times 32 \times 3$ Cifar-10	MLP	3	256	–	–	4096	–
$\text{C10T}_{(3,512)}^{\text{MLP}}$	$32 \times 32 \times 3$ Cifar-10	MLP	3	512	–	–	4096	–
$\text{C10T}_{(3,1024)}^{\text{MLP}}$	$32 \times 32 \times 3$ Cifar-10	MLP	3	1024	–	–	4096	–
$\text{C10T}_{(3,2048)}^{\text{MLP}}$	$32 \times 32 \times 3$ Cifar-10	MLP	3	2048	–	–	4096	–
$\text{C10T}_{(3,4096)}^{\text{MLP}}$	$32 \times 32 \times 3$ Cifar-10	MLP	3	4096	–	–	4096	–
$\text{C10T}_{(3,8192)}^{\text{MLP}}$	$32 \times 32 \times 3$ Cifar-10	MLP	3	8192	–	–	4096	–
$\mathcal{T}_{(3,192)}^{\text{ViT}}$	$32 \times 32 \times 3$ ImageNet	ViT	3	192	3	768	4096	–
$\mathcal{T}_{(3,384)}^{\text{ViT}}$	$32 \times 32 \times 3$ ImageNet	ViT	3	384	6	1536	4096	–
$\mathcal{T}_{(3,768)}^{\text{ViT}}$	$32 \times 32 \times 3$ ImageNet	ViT	3	768	8	3072	4096	–
$\mathcal{T}_{(3,1024)}^{\text{ViT}}$	$32 \times 32 \times 3$ ImageNet	ViT	3	1024	8	4096	4096	–
$\mathcal{T}_{(3,2048)}^{\text{ViT}}$	$32 \times 32 \times 3$ ImageNet	ViT	3	2048	16	8192	4096	–
$\mathcal{T}_{(3,3072)}^{\text{ViT}}$	$32 \times 32 \times 3$ ImageNet	ViT	3	3072	16	12288	4096	–
$\mathcal{T}_{(3,192)}^{\text{LM}}$	LM1B, 32k Vocab	Transformer LM	3	192	3	768	4096	64
$\mathcal{T}_{(3,384)}^{\text{LM}}$	LM1B, 32k Vocab	Transformer LM	3	384	6	1536	4096	64
$\mathcal{T}_{(3,768)}^{\text{LM}}$	LM1B, 32k Vocab	Transformer LM	3	768	8	3072	4096	64
$\mathcal{T}_{(3,1024)}^{\text{LM}}$	LM1B, 32k Vocab	Transformer LM	3	1024	8	4096	4096	64
$\mathcal{T}_{(3,2048)}^{\text{LM}}$	LM1B, 32k Vocab	Transformer LM	3	2048	16	8192	4096	64
$\mathcal{T}_{(3,3072)}^{\text{LM}}$	LM1B, 32k Vocab	Transformer LM	3	3072	16	12288	4096	64
$\mathcal{DT}_{(16,512)}^{\text{MLP}}$	32×32 ImageNet	MLP	16	512	–	–	4096	–
$\mathcal{DT}_{(16,192)}^{\text{ViT}}$	32×32 ImageNet	ViT	16	192	3	768	4096	–
$\mathcal{DT}_{(16,192)}^{\text{LM}}$	LM1B	Transformer LM	16	192	3	768	4096	–

Appendix F. Extended Results

F.1. Evaluating pre-activation stability

We now verify that desiderata J.1 of [36] is satisfied empirically. In Figure 9, we report the evolution of the coordinate-wise standard deviation of the difference between initial ($t=0$) and current (t) second-layer pre-activations of an MLP during the first 500 steps of training for a single trial. We observe that all models parameterized in μ P enjoy stable coordinates across widths, while the pre-activations of the larger models in SP blow up after a number of training steps. Notably, SP Adam’s pre-activations blow up immediately, while LO_S and LO_M take longer to blow up and have a more erratic pattern; we hypothesize that this is a side effect of meta-training where the optimizers may learn to keep pre-activations small by rescaling updates. Section J of the appendix contains similar plots for the remaining layers of the MLP which show similar trends. In summary, we find, empirically, that pre-activations of μ LOs and μ Adam are similarly stable across widths, while the activations of SP Adam and SP LOs both blow up but behave qualitatively differently.

Appendix G. Hand Designed Optimizer Hyperparameter Tuning

G.1. Tuning μ Adam

We tune the μ Adam baseline on the largest meta-training seen by our learned optimizers. μ Adam_M was, therefore, tuned using a 1024 width MLP for $32 \times 32 \times 3$ ImageNet classification. As mentioned in section 3, we tune the learning rate and three multipliers: input multiplier, output multiplier, and the hidden learning rate multiplier. These multipliers correspond to adding a tunable constant to the pre-activation multiplier for input weights, the pre-activation multiplier for output weights, and the Adam LR for hidden weights (e.g., in Table 8 of [36]). Specifically, we search for the learning rate in $\{0.1, 0.01, 0.001, 0.0001\}$ and for each multiplier in $\{2^{-4}, 2^{-2}, 1, 2^2, 2^4\}$. This results in a grid search of 500 configurations, whose optimal values are reported in table 3.

Table 3: **Best hyperparameters values for μ Adam baseline.** μ Adam is tuned to optimize 3-layer $W=1024$ MLP for $32 \times 32 \times 3$ ImageNet classification.

Optimizer	LR	Input Multiplier	Output Multiplier	Hidden LR Multiplier
μ Adam	0.1	0.25	0.25	4

G.2. Tuning AdamW

We tune the AdamW baseline on the largest meta-training seen by our learned optimizers. AdamW was, therefore, tuned using a 1024 width MLP for $32 \times 32 \times 3$ ImageNet classification. As mentioned in section 3, we tune the learning rate, betas, and weight decay: LR, β_1 , β_2 , and the weights decay. Specifically, we search over the values of each hyperparameter reported in Table 4. This results in a grid search of 500 configurations, whose optimal values are reported in table 5.

Table 4: **Grid search values used for AdamW**. Similar to the μ Adam baseline, we tune all optimizers on a 3-layer $W=1024$ MLP ImageNet classification task and use a budget of approximately 500 total runs. We tune LR, β_1 , β_2 , and weight decay to minimize training loss after 1000 steps.

Optimizer	LR	β_1	β_2	weight decay	Total runs
SP AdamW	Log Sample 14 from $[10^{-5}, 0.1]$	{0.9,0.95,0.99}	{0.95,0.99,0.999}	{0.1,0.01,0.001,0.0001}	504

Table 5: **Optimal Hyperparameters Found AdamW**. Similar to the μ Adam baseline, we tune all optimizers on a 3-layer $W=1024$ MLP ImageNet classification task and use a budget of approximately 500 total runs.

Optimizer	LR	β_1	β_2	weight decay	Total runs
SP AdamW	0.000702	0.9	0.95	0.0001	504

Appendix H. Meta-training with μ LOs

General meta-training setup Each learned optimizer is meta-trained for 5000 steps of gradient descent using AdamW [16] and a linear warmup and cosine annealing schedule. We use PES [29] to estimate meta-gradients with a truncation length of 50 steps and sampling 8 perturbations at each step with standard deviation 0.01. For the inner optimization task, we used a maximum unroll length of 1000 iterations; that is, all our learned optimizers see at most 1000 steps of the inner optimization problem during meta-training. Unlike with μ Adam, we do not tune the μ P multipliers when meta-training μ LO_S and μ LO_M, instead, we set the all to 1. All optimizers are meta-trained on a single A6000 GPU. μ LO_S and LO_S take 8 hours each to meta-train, while μ LO_M and LO_M take 103 hours.

Ablating meta-training hyperparameters for μ LOs While there are very few differences between μ LOs and SP LOs, the effective step size for hidden layers is changed (see eq. 1) which could alter the optimal meta-training hyperparameters. Consequently, we conduct an ablation study on hyperparameter choices for μ LO_S. Specifically, using AdamW and gradient clipping with a linear warmup and cosine annealing LR schedule, we meta-train μ LO_S to train 3-layer width 128 MLPs to classify $32 \times 32 \times 3$ ImageNet Images. By default, we warmup linearly for 100 steps to a maximum learning rate of $3e-3$ and anneal the learning rate for 4900 steps to a value of $1e-3$ with $\lambda_1 = 0.001$ (from equation 1) and sampling 8 perturbations per step in PES[29]. The above ablation varies the maximum learning rate $\in \{1e-2, 3e-3, 1e-3\}$ (always using 100 steps of warmup and decaying to $0.3 \times \text{MaxLR}$), $\lambda_1 \in \{0.001, 0.01, 0.1\}$, the number of steps (5k or 10k), and the number of perturbations (8 or 16). We observe that using all default values except for $\lambda_1 = 0.01$ yields one of the best solutions while being fast to train and stable during meta-training. We, therefore, select these hyperparameters to meta-train μ LO_S and μ LO_M.

μ P at Meta-training time While we use the same μ P at meta-training and testing time, it is important to consider meta-training tasks that have similar training trajectories to their infinite width counterparts. In [36], authors provide discussions of these points for zero-shot hyperparameter transfer. Two notable guidelines are to initialize the output weight matrix to zero (as it will approach zero in the limit) and to use a relatively large key size when meta-training transformers. For all our tasks, we initialize the network’s final layer to zeros. While we do not meta-train on transformers, we suspect that the aforementioned transformer-specific guidelines may be useful.

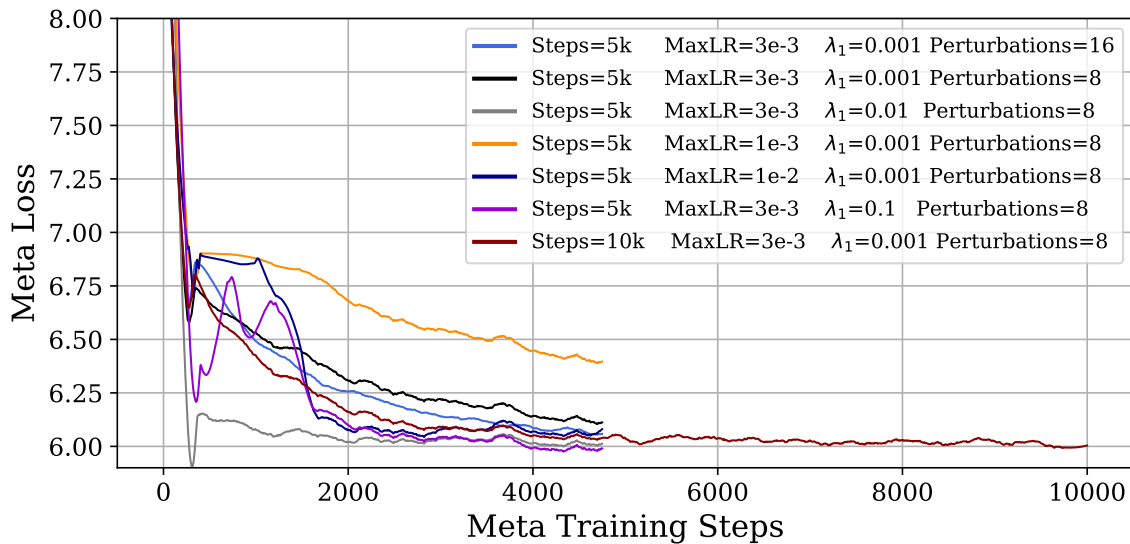


Figure 6: **Ablating Meta-training Hyperparameter for μLO_S .** Using AdamW with a linear warmup and cosine annealing schedule, we meta-train μLO_S to train 3-layer width 128 MLPs for classifying $32 \times 32 \times 3$ ImageNet Images. By default, we warmup linearly for 100 steps to a maximum learning rate of $3e-3$ and anneal the learning rate for 4900 steps to a value of $1e-3$ with $\lambda_1 = 0.001$ (from equation 1) and sampling 8 perturbations per step in PES[29]. The above ablation varies the maximum learning rate $\in \{1e-2, 3e-3, 1e-3\}$ (always using 100 steps of warmup and decaying to $0.3 \times \text{MaxLR}$), $\lambda_1 \in \{0.001, 0.01, 0.1\}$, the number of steps (5k or 10k), and the number of perturbations (8 or 16). We observe that using all default values except for $\lambda_1 = 0.01$ yields one of the best solutions while being fast to train and stable during meta-training.

Appendix I. Features of the learned optimizer

Table 6: **Per-parameter features used as input to our learned optimizers.** All the coefficients, β_i , are learnable parameters adjusted during meta-optimization. We replicate the table of [12] for convenience.

Description	value
parameter value	w_t
3 momentum values with coefficients $\beta_1, \beta_2, \beta_3$	$m_{t,i} = \beta_i m_{t-1,i} + (1 - \beta_i) g_t$
second moment value computed from g_t with decay β_4	$v_t = \beta_4 v_{t-1} + (1 - \beta_4) g_t^2$
3 values consisting of the three momentum values normalized by the square root of the second moment	$\frac{m_{t,i}}{\sqrt{v}}$
the reciprocal square root of the second moment value	$\frac{1}{\sqrt{v}}$
3 Δ_t Adafactor normalized values	$g_t \times \text{ROW FACTOR} \times \text{COLUMN FACTOR}$
3 tiled Adafactor row features with coefficients $\beta_5, \beta_6, \beta_7$, computed from g_t	$r_{t,i} = \beta_i r_{t-1,i} + (1 - \beta_i) \text{ROW_MEAN}(\Delta_t^2)$
3 tiled Adafactor column feature with coefficients $\beta_5, \beta_6, \beta_7$ computed from g_t	$c_{t,i} = \beta_i c_{t-1,i} + (1 - \beta_i) \text{COL_MEAN}(\Delta_t^2)$
the reciprocal square root of the previous 6 features	$\frac{1}{\sqrt{r_{t,i} \text{ OR } c_{t,i}}}$
3 m Adafactor normalized values	$m_{t,i} \times \text{ROW FACTOR} \times \text{COLUMN FACTOR}$

Appendix J. Coordinate evolution of MLP layers in μ P for different optimizers

The following section presents the continuation of our experiments comparing pre-activation growth during training for SP LOs and μ LOs with different meta-training recipes, SP adam, and μ Adam.

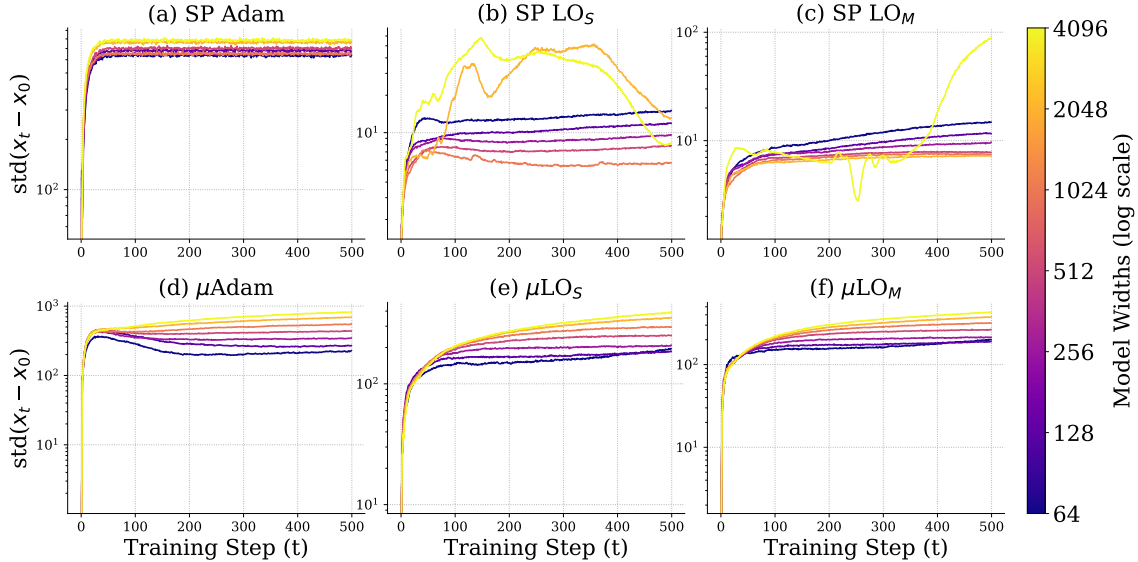


Figure 7: **Layer 0 pre-activations behave harmoniously in μ P for LOs and Adam alike.** We report the evolution of coordinate-wise standard deviation between the difference of initial and current second-layer pre-activations. We observe that all models parameterized in μ P enjoy stable coordinates across widths, while the pre-activations of larger-width models in SP blow up after a number of training steps. All plots report these metrics for the first 500 steps of a single training run.

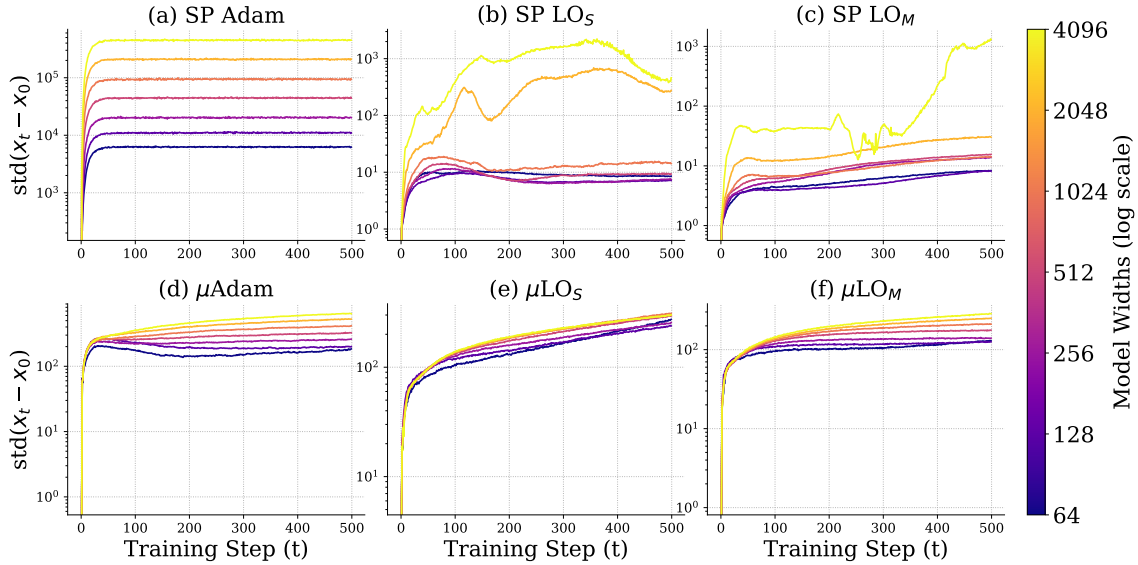


Figure 8: **Layer 1 pre-activations behave harmoniously in μ P for LOs and Adam alike.** We report the evolution of coordinate-wise standard deviation between the difference of initial and current second-layer pre-activations. We observe that all models parameterized in μ P enjoy stable coordinates across widths, while the pre-activations of larger-width models in SP blow up after a number of training steps. All plots report these metrics for the first 500 steps of a single training run.

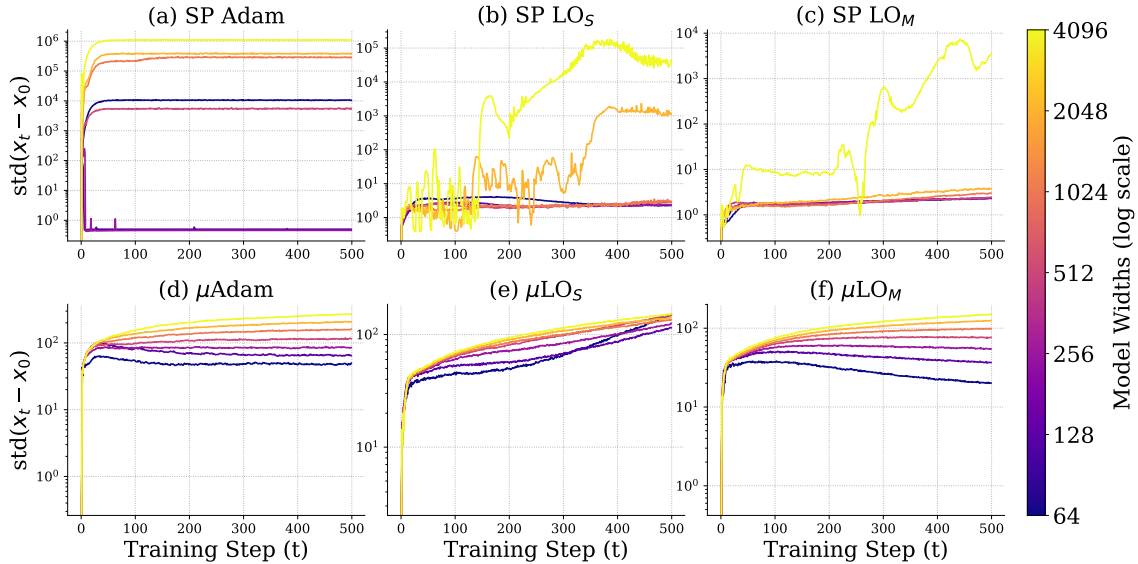


Figure 9: **Layer 2 pre-activations behave harmoniously in μ P for μ LOs and μ Adam alike.** We report the evolution of coordinate-wise standard deviation of the difference between the initial ($t = 0$) and t -th second-layer pre-activations of an MLP during training for the first 500 steps of a single run (the remaining layers behave similarly, see Sec. J). We observe that all models parameterized in μ P enjoy stable coordinates across widths, while the pre-activations of larger-width models in SP blow up after a number of training steps.

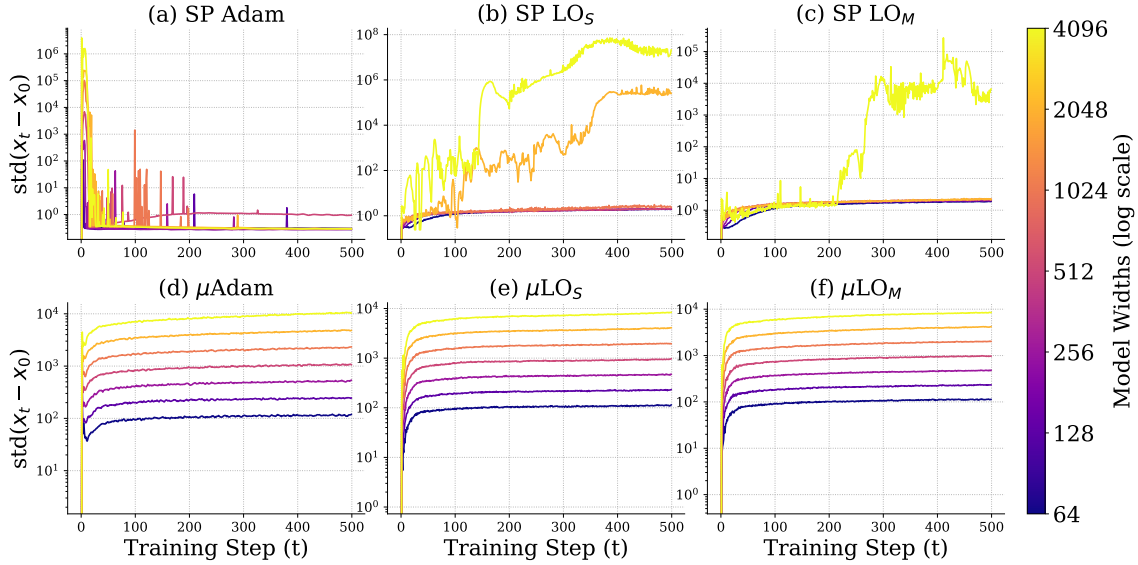


Figure 10: **Layer 3 pre-activations behave harmoniously in μ P for LOs and Adam alike.** We report the evolution of coordinate-wise standard deviation between the difference of initial and current second-layer pre-activations. We observe that all models parameterized in μ P enjoy stable coordinates across widths, while the pre-activations of larger-width models in SP blow up after a number of training steps. All plots report these metrics for the first 500 steps of a single training run.

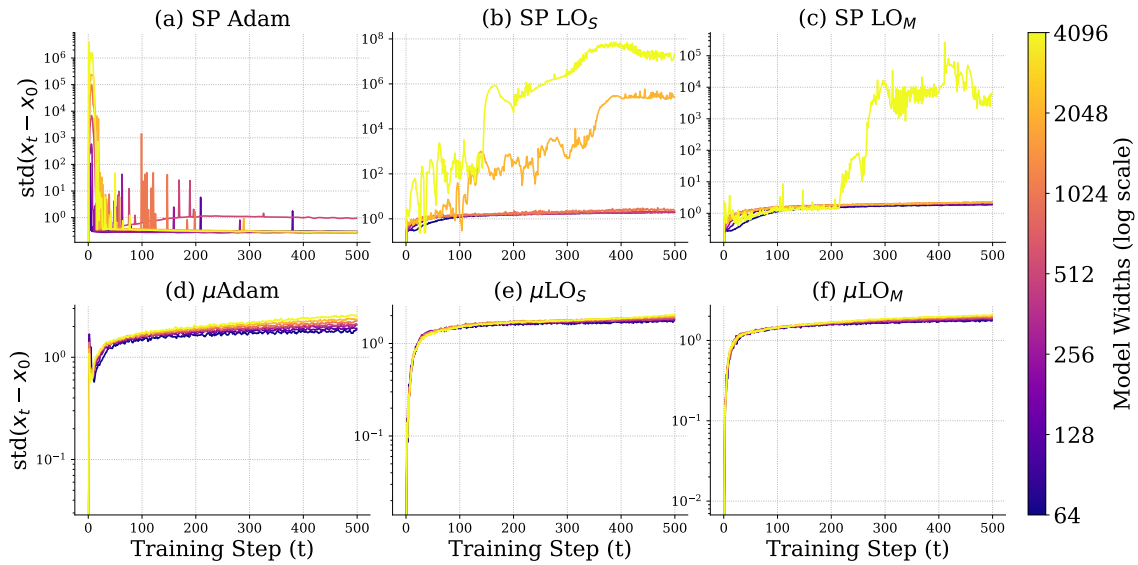


Figure 11: **Logits behave harmoniously in μ P for LOs and Adam alike.** We report the evolution of coordinate-wise standard deviation between the difference of initial and current second-layer pre-activations. We observe that all models parameterized in μ P enjoy stable logits across widths, while the pre-activations of larger-width models in SP blow up after a number of training steps. All plots report these metrics for the first 500 steps of a single training run.