

# $\mu$ LO: Compute-Efficient Meta-Generalization of Learned Optimizers

by **Benjamin Thérien, Charles-Étienne Joseph, Boris Knyazev, Edouard Oyallon, Irina Rish, Eugene Belilovsky**

ChatGPT

“Here is an image of a learned optimizer training a gigantic neural network in a futuristic setting.”

# Overview

**The case for learning our  
optimizers  
+  
Background**

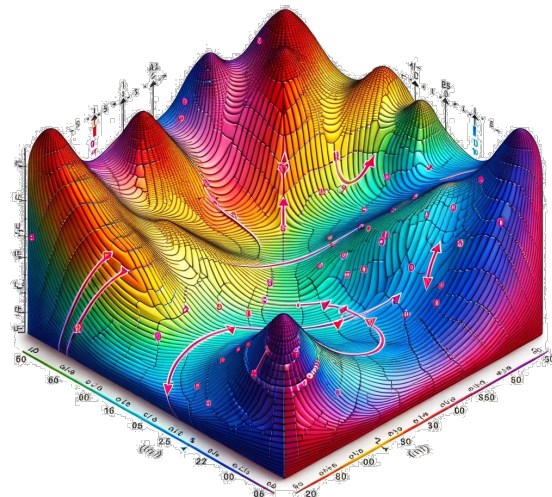
**$\mu$ LOs: results,  
performance, and  
outlook**

**A case study of  
VeLO  
(where we are  
now)**

**The need for Maximal  
Update  
Parameterization( $\mu$ P)  
for Learned  
Optimizers**

# Why should we learn optimizers?

- Gradient-based optimizers such as Adam, SGD, etc... are prevalent in DL.
- But, DL optimization problems are highly non-convex and theory is hard.
- No reason to expect our optimizers will converge to the global optimum at the optimal rate.



**Can we meta-learn better optimization algorithms?**

## The Bitter Lesson:

**“General methods that leverage computation are ultimately the most effective, and by a large margin.”**

- Rich Sutton, 2019

# Learned Optimizer Background

# General Optimizer Update

**Optimizee  
Parameters**

$$\boldsymbol{\theta}_t = f_{\phi}(\boldsymbol{\theta}_{t-1}, \mathbf{U}_{t-1})$$

**Optimizer  
Parameters**

**Optimizer  
Input**

# Learned Optimizer Update

**Optimizee  
Parameters**

$$\boldsymbol{\theta}_t = f_\phi(\boldsymbol{\theta}_{t-1}, \mathbf{U}_{t-1}) = \boldsymbol{\theta}_{t-1} - \frac{\lambda_1 \mathbf{d}_\phi e^{(\lambda_2 \mathbf{m}_\phi)}}{\text{LO Update}}$$

$$g_\phi(\boldsymbol{\theta}_{t-1}, \mathbf{U}_{t-1}) = [\mathbf{d}_\phi, \mathbf{m}_\phi]$$

**Optimizer  
Neural Network**

# Learned Optimizer Input Features

Type	#	Description	Equation
Accumulators	3	Momentum accumulators with coefficients $\beta_i, i \in \{1, 2, 3\}$ .	$m_{t,i} = \beta_i m_{t-1,i} + (1 - \beta_i) \nabla_t$
	1	Second moment accumulator with coefficients $\beta_4$ .	$v_t = \beta_4 v_{t-1} + (1 - \beta_4) \nabla_t^2$
	3	Adafactor row accumulator with coefficients $\beta_i, i \in \{5, 6, 7\}$ .	$r_{t,i} = \beta_i r_{t-1,i} + (1 - \beta_i) \text{row\_mean}(\nabla_t^2)$
	3	Adafactor accumulator with coefficients $\beta_i, i \in \{5, 6, 7\}$ .	$c_{t,i} = \beta_i c_{t-1,i} + (1 - \beta_i) \text{col\_mean}(\nabla_t^2)$

# Learned Optimizer Input Features

Type	#	Description	Equation
Accumulators	3	Momentum accumulators with coefficients $\beta_i, i \in \{1, 2, 3\}$ .	$m_{t,i} = \beta_i m_{t-1,i} + (1 - \beta_i) \nabla_t$
	1	Second moment accumulator with coefficients $\beta_4$ .	$v_t = \beta_4 v_{t-1} + (1 - \beta_4) \nabla_t^2$
	3	Adafactor row accumulator with coefficients $\beta_i, i \in \{5, 6, 7\}$ .	$r_{t,i} = \beta_i r_{t-1,i} + (1 - \beta_i) \text{row\_mean}(\nabla_t^2)$
	3	Adafactor accumulator with coefficients $\beta_i, i \in \{5, 6, 7\}$ .	$c_{t,i} = \beta_i c_{t-1,i} + (1 - \beta_i) \text{col\_mean}(\nabla_t^2)$

# Learned Optimizer Input Features

Type	#	Description	Equation
Accumulators	3	Momentum accumulators with coefficients $\beta_i, i \in \{1, 2, 3\}$ .	$m_{t,i} = \beta_i m_{t-1,i} + (1 - \beta_i) \nabla_t$
	1	Second moment accumulator with coefficients $\beta_4$ .	$v_t = \beta_4 v_{t-1} + (1 - \beta_4) \nabla_t^2$
	3	Adafactor row accumulator with coefficients $\beta_i, i \in \{5, 6, 7\}$ .	$r_{t,i} = \beta_i r_{t-1,i} + (1 - \beta_i) \text{row\_mean}(\nabla_t^2)$
	3	Adafactor accumulator with coefficients $\beta_i, i \in \{5, 6, 7\}$ .	$c_{t,i} = \beta_i c_{t-1,i} + (1 - \beta_i) \text{col\_mean}(\nabla_t^2)$

# Learned Optimizer Input Features

Type	#	Description	Equation
<b>Accumulator Features</b>	3	Momentum values normalized by the square root of the second moment for $i \in \{5, 6, 7\}$ .	$\frac{m_{t,i}}{\sqrt{v}}$
	1	The reciprocal square root of the second moment value.	$\frac{1}{\sqrt{v}}$
	6	The reciprocal square root of the Adafactor accumulators.	$\frac{1}{\sqrt{r_{t,i} \text{ OR } c_{t,i}}}$
	3	Adafactor gradient features for $i \in \{5, 6, 7\}$ .	$\nabla_t \times r_{t,i} \times c_{t,i}$
	3	Adafactor momentum features for $i, j \in \{(5, 1), (6, 2), (7, 3)\}$ .	$m_{t,j} \times r_{t,i} \times c_{t,i}$
<b>Time Features</b>	11	Time Features for $x \in \{1, 3, 10, 30, 100, 300, 1000, 3000, 10^4, 3 \cdot 10^4, 10^5\}$ .	$\tanh\left(\frac{t}{x}\right)$
<b>Parameters</b>	1	Parameter value.	$w_t$

# Learned Optimizer's Objective

4.0 Learned Optimization

$$\min_{\phi} \mathbb{E}_{(\mathcal{D}, \mathcal{L}, \theta_0) \sim \mathcal{T}} \left[ \mathbb{E}_{(X, Y) \sim \mathcal{D}} \left[ \frac{1}{T} \sum_{t=1}^T \mathcal{L}(X, Y; f_{\phi}(\theta_{t-1}, \mathbf{U}_{t-1})) \right] \right]$$

Diagram labels and connections:

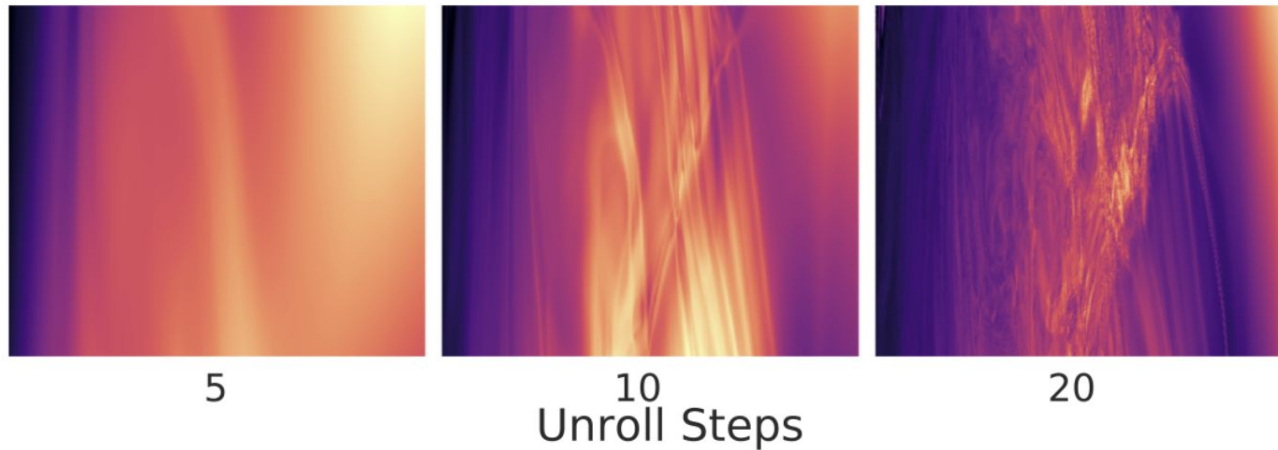
- Learned Optimizer** points to  $\phi$ .
- Optimizee initialization** points to  $\theta_0$ .
- Data distribution (e.g., web text, ImageNet)** points to  $\mathcal{D}$ .
- Batch of training data** points to  $(X, Y)$ .
- unroll steps** points to  $T$ .
- Optimizee step t Loss** points to  $\mathcal{L}(X, Y; f_{\phi}(\theta_{t-1}, \mathbf{U}_{t-1}))$ .
- Learned optimizer** points to  $f_{\phi}$ .

- 1) Sample a data distribution  $\mathbf{D}$ , Loss  $\mathbf{L}$ , and weight  $\theta$ .
- 2) Sample  $\mathbf{X}, \mathbf{Y}$  from  $\mathbf{D}$ .
- 3) Compute the loss,  $\mathbf{L}$ , for the current optimizee weights  $\theta$ .
- 4) Use  $\phi$  to update  $\mathbf{w}$ .
- 5) Repeat for  $\mathbf{T}$  steps.

# Noisy Meta-loss Landscapes



(Metz et al., 2019)



- Unrolled optimization problems have many steps
- Longer unrolls = Noisy gradients
- We should smooth the meta loss

# Zeroth Order Meta-gradient Estimation

$$\nabla_{\phi}^{\text{ES-A}} = \frac{1}{N\sigma^2} \sum_{i=1}^{N/2} \epsilon^{(i)} \left( L(\phi + \epsilon^{(i)}) - L(\phi - \epsilon^{(i)}) \right)$$

- By sampling epsilon from a gaussian we smooth the loss
- Zeroth order methods are used instead of backprop
- Antithetic sampling is used to reduce variance

# **VeLO Optimizer: A Case Study (where we are now)**

# Scaling learned optimization leads to faster per-step optimization

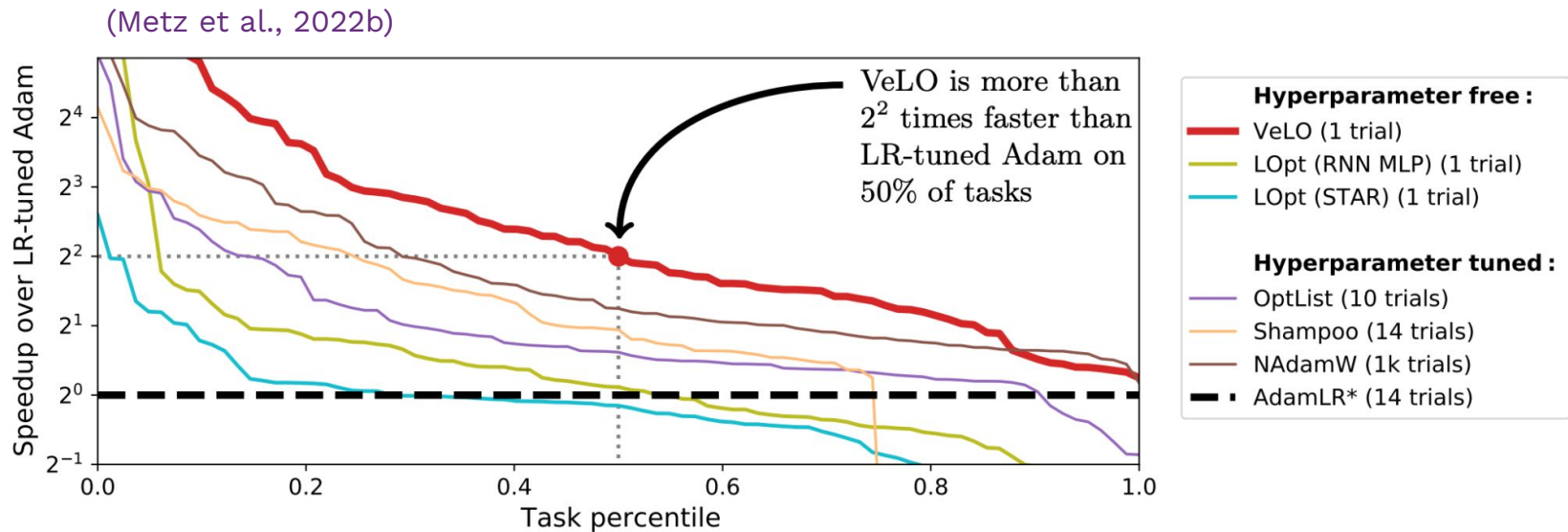
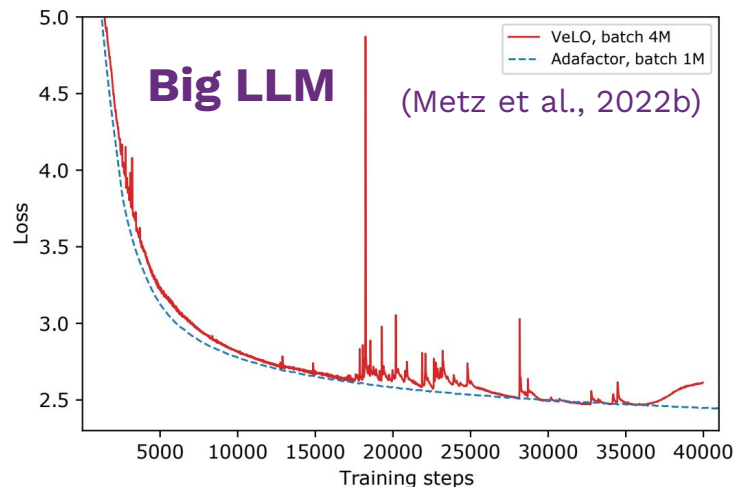
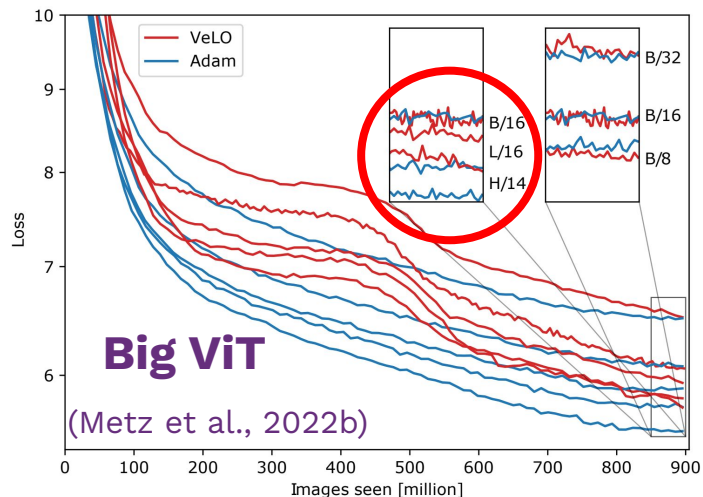


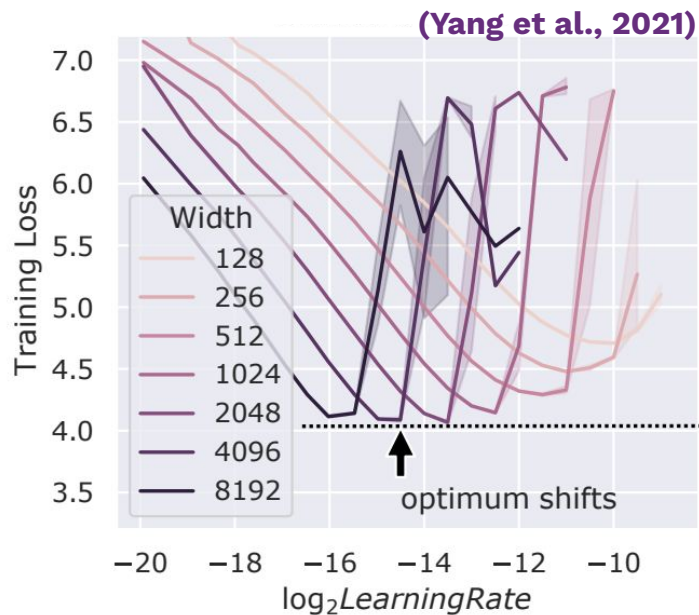
Figure 1: **Optimizer performance on the 83 canonical tasks in the VeLOdrome benchmark.**

# VeLO is outperformed by hand-designed optimizers for the largest models



- VeLO is tested on big models
- VeLO fails
- Why? VeLO is not trained on big models

# The optimal learning rate depends on width



- Optimal LR depends on width in Standard Parameterization

# $\mu$ P Background; Parameterization

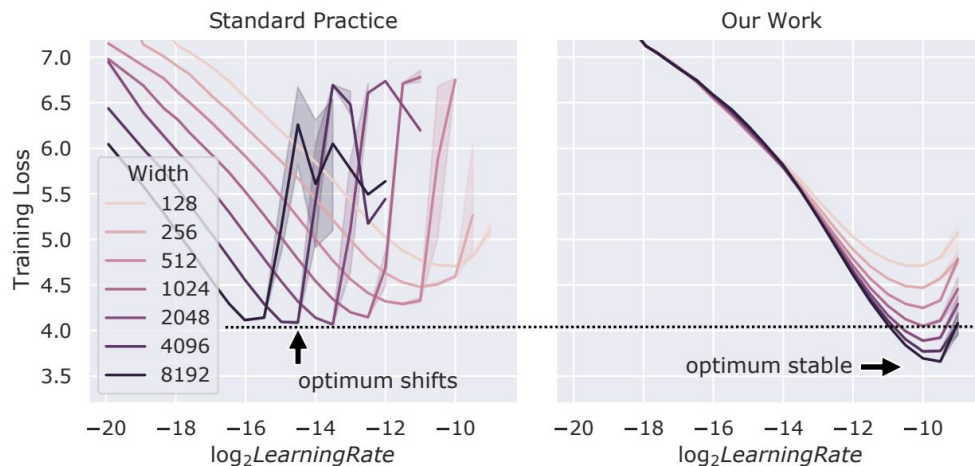
(Yang et al., 2021)

	Input weights & all biases		Output weights		Hidden weights	
Init. Var.		$1/\text{fan\_in}$	<b>1</b>	$(1/\text{fan\_in})$		$1/\text{fan\_in}$
Multiplier		<b>1</b>	$1/\text{fan\_in}$	(1)		<b>1</b>
SGD LR	<b>fan_out</b>	(1)	<b>fan_in</b>	(1)		<b>1</b>
Adam LR		<b>1</b>		<b>1</b>	$1/\text{fan\_in}$	(1)

- A parameterization specifies:
  - initialization variance
  - pre-activation multipliers
  - LR multipliers

# $\mu$ P enables HP transfer from small width to large width models

(Yang et al., 2021)



- (**LEFT**) Optimal LR depends on width in SP
- (**RIGHT**) Optimal LR is the same for all widths in  $\mu$ P
- $\mu$ P stabilizes training for infinite width

# $\mu$ P for learned optimizers

# How $\mu P$ Impacts the LO Update

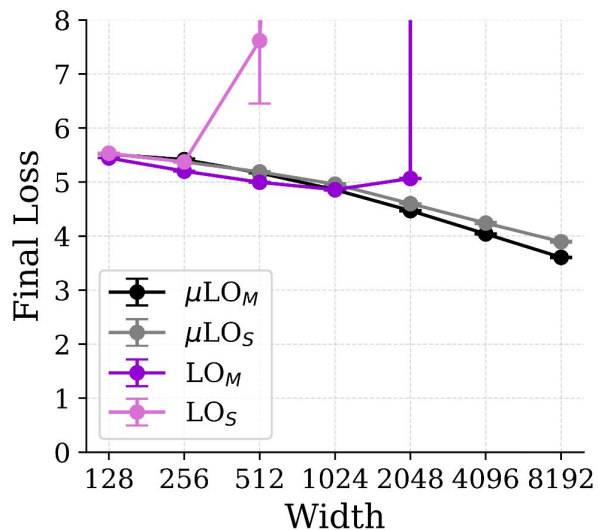
**Updated Parameter Value**       **$\mu P$  LR Multiplier**      **Direction**

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

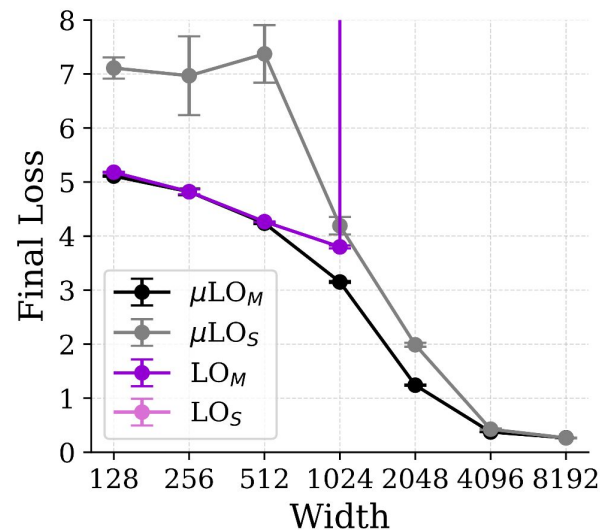
**Magnitude**

- LO  $\mu P$  requires scaling the LR for hidden weights

# $\mu$ LO Meta-training Recipe



(a) Iteration 1000

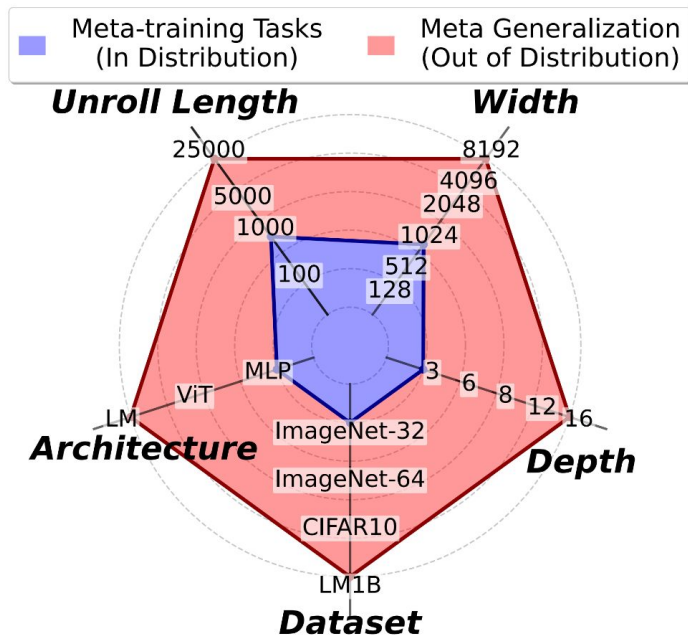


(b) Iteration 5000

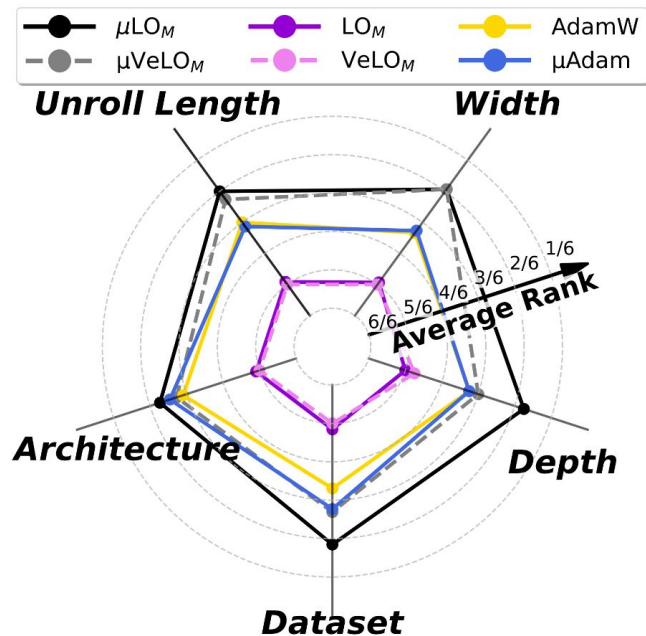
- **3-task Meta-training** v.s. **1-task Meta-training**
- **3-task Meta-training** consistently performs better
- We use **3-task Meta-training** in our experiments

# Results and Outlook of $\mu$ LOs

# Results summary



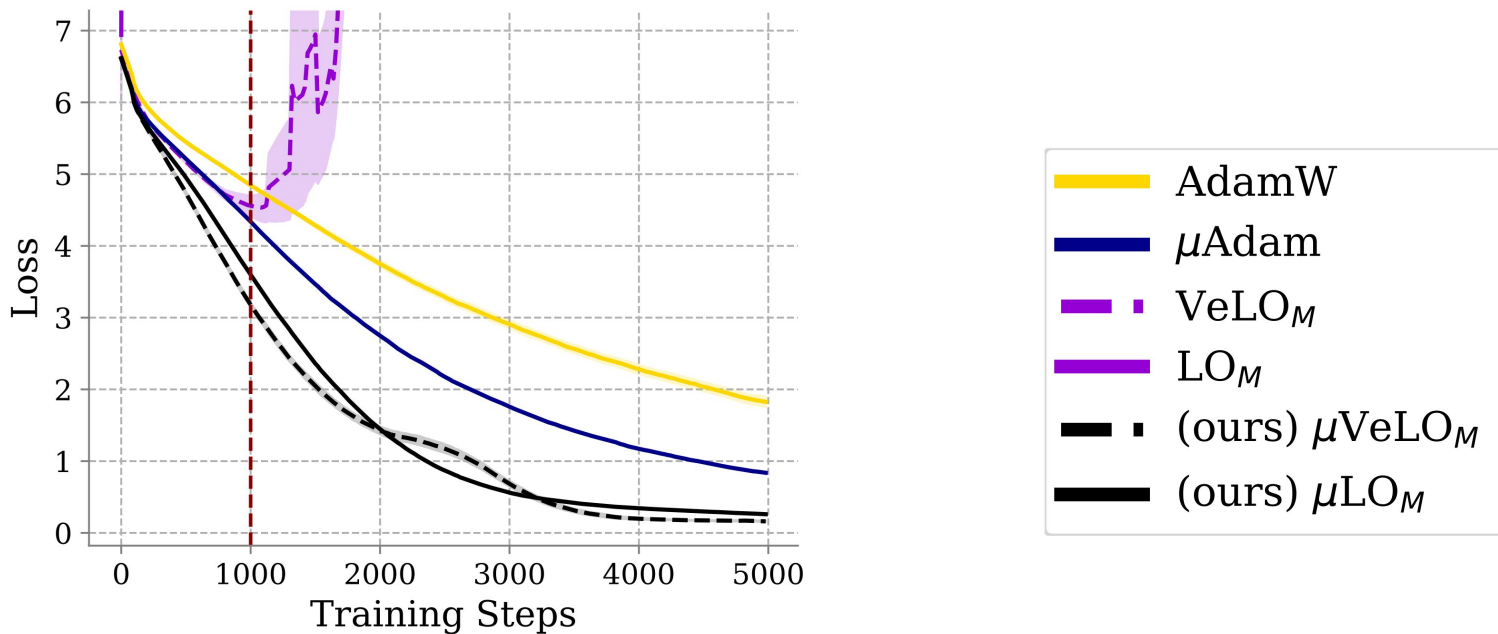
(a) Axes of Meta-Generalization



(b) Performance by Average Rank

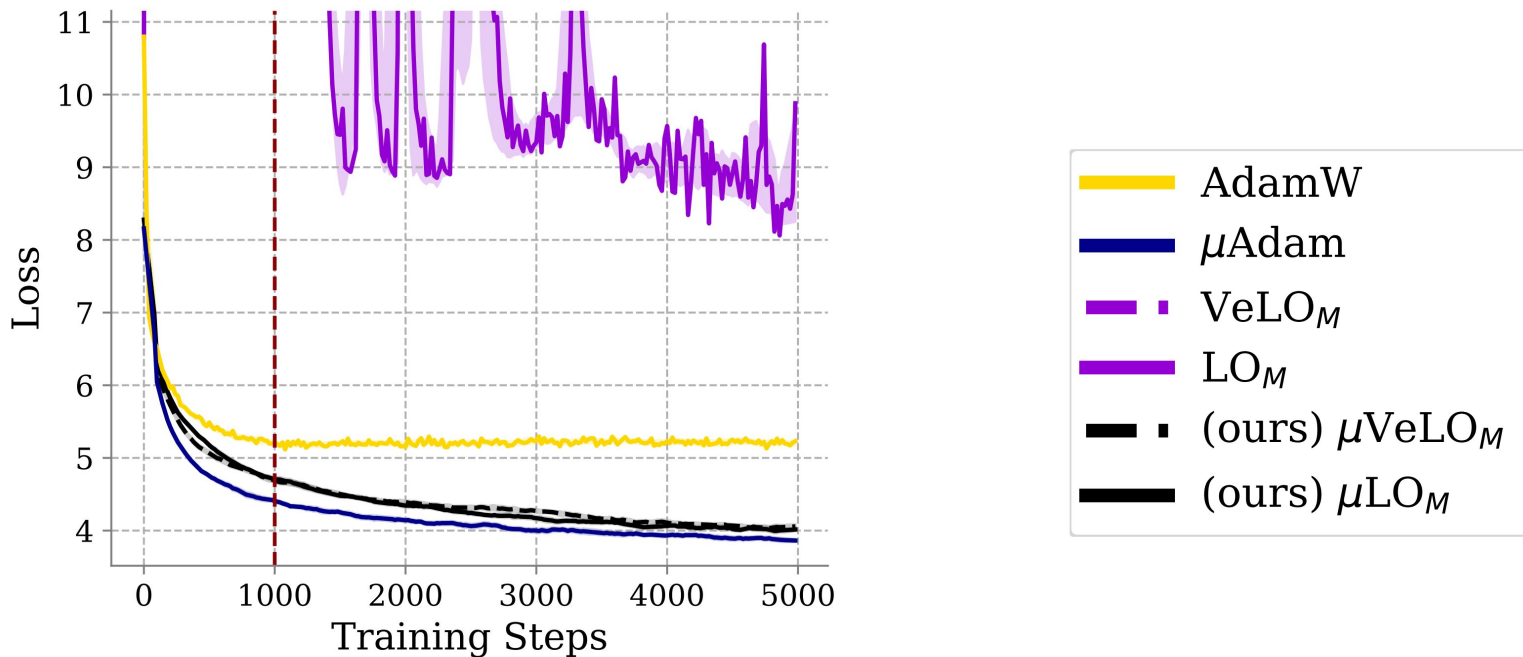
- $\mu LO$ s reach the highest average rank of all tasks

# Width Generalization - Imagenet 32x32x3



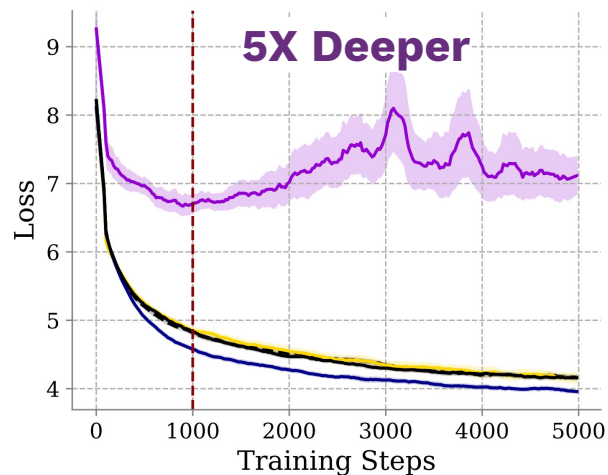
(a) MLP IN32 W=8192

# Width Generalization - Transformer LM

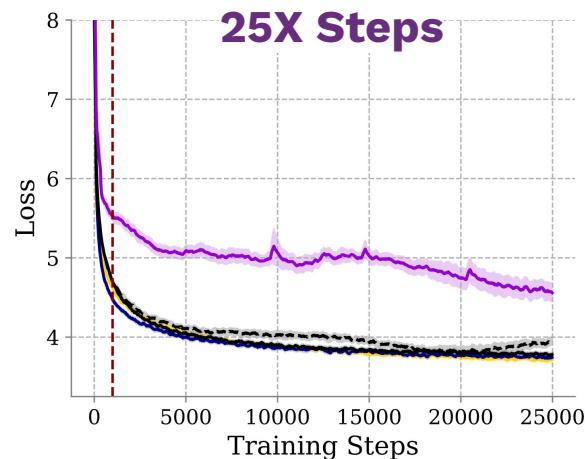


(d) LM  $W=4096$

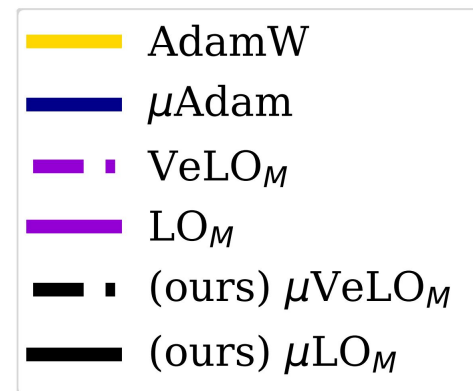
# Deeper networks & Longer training



(b) LM W=1024 D=16

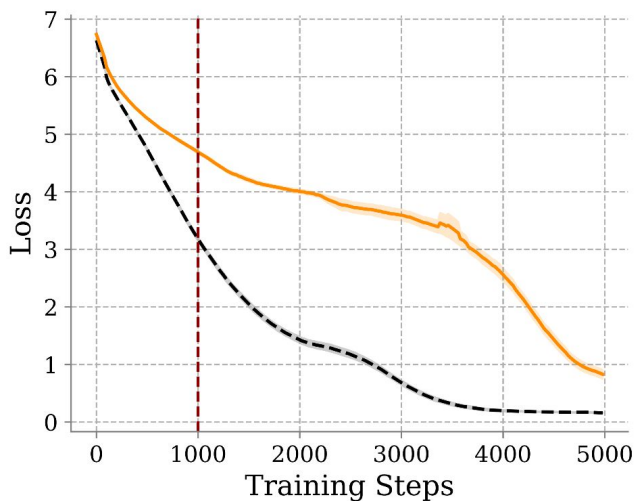


(b) LM W=1024

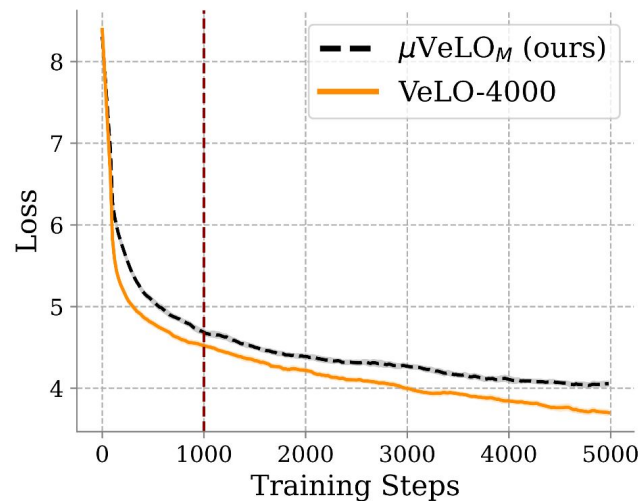


- We test on deeper networks and for longer training.
- We observe empirically, that  $\mu$ LOs improve generalization in both these cases.

# Comparison with VeLO



(a) MLP IN32 W=8192



(b) LM W=4096

- We outperform VeLO when testing on the meta-training distribution.
- VeLO outperforms on OOD tasks.
- Will expanding  $\mu\text{LOs}$  meta-training distributions bridge the gap?

# Conclusion: Learned optimizers should be meta-trained in $\mu$ -parameterization.

- We derived  $\mu$ -parameterization for popular learned optimizers (VeLO and `small_fc_lopt`) and proposed a meta-training recipe for  $\mu$ LOs.
- $\mu$ LOs generalize significantly better to wider networks than SP LOs and outperform VeLO on wider in-distribution tasks!
- $\mu$ LOs also improve generalization to deeper networks (5X meta-training) and longer training horizons (25X meta-training).