

# Warmstarting for Scaling Language Models

Neeratyoy Mallik\*, Maciej Janowski\*, Johannes Hog, Herilalaina Rakotoarison, Aaron Klein, Josif Grabocka, Frank Hutter

\* Equal contribution

Improving convergence rate of larger models by warmstarting training from a smaller model under Chinchilla compute-optimal training.

- Scaling studies and large model trainings often do not have directly transferable hyperparameter settings.
- Prohibitive tuning cost at large scales lead to custom tuning decisions across different setups.
- Reducing training cost at larger scales by reusing trained smaller models.
- Leverage tuning decisions made at smaller model scales.

Scaling best found learning rate at smallest model scale using  $\mu P$ .

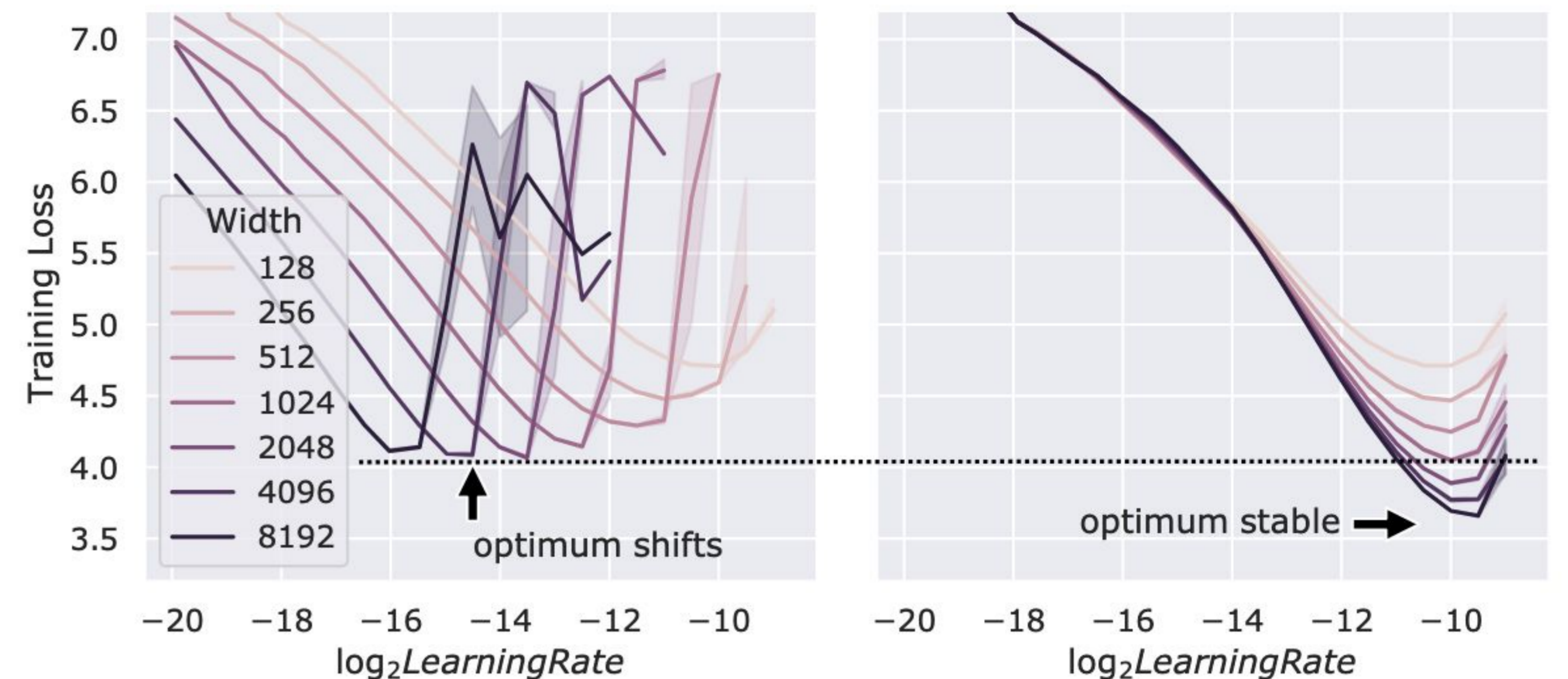


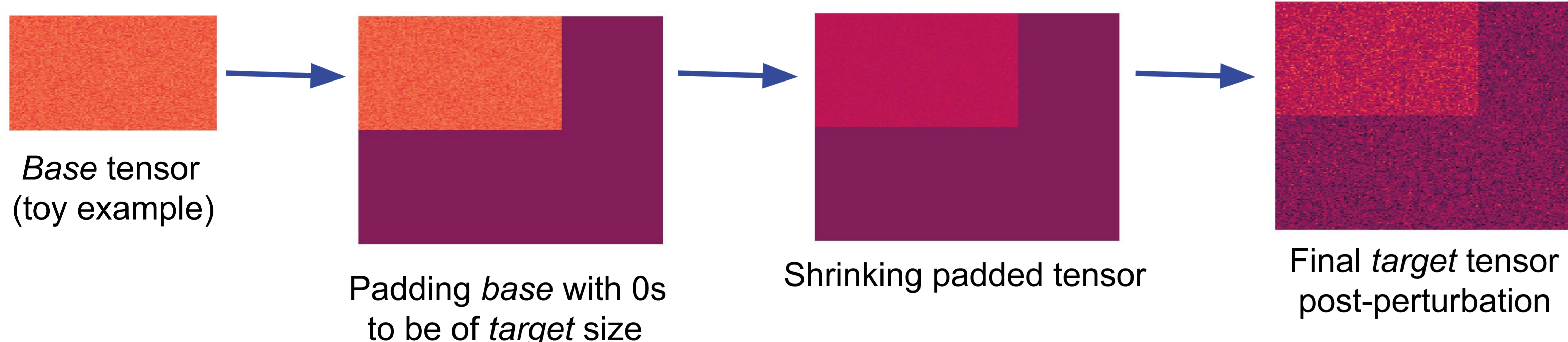
Figure source: [1]

## Our Approach

Initialize the larger model training run as a scaled up continual learning over the smaller model, assuming  $\mu P$ -enabled training pipeline.

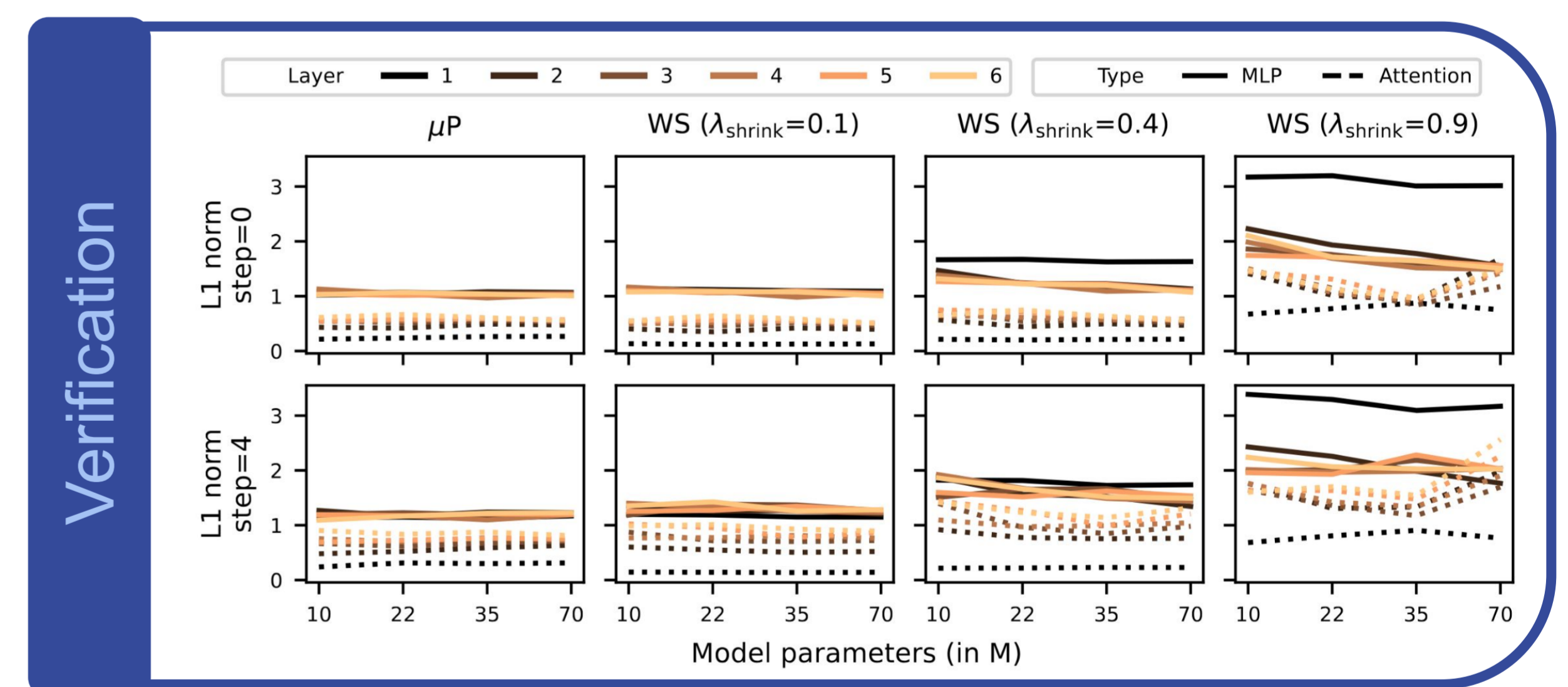
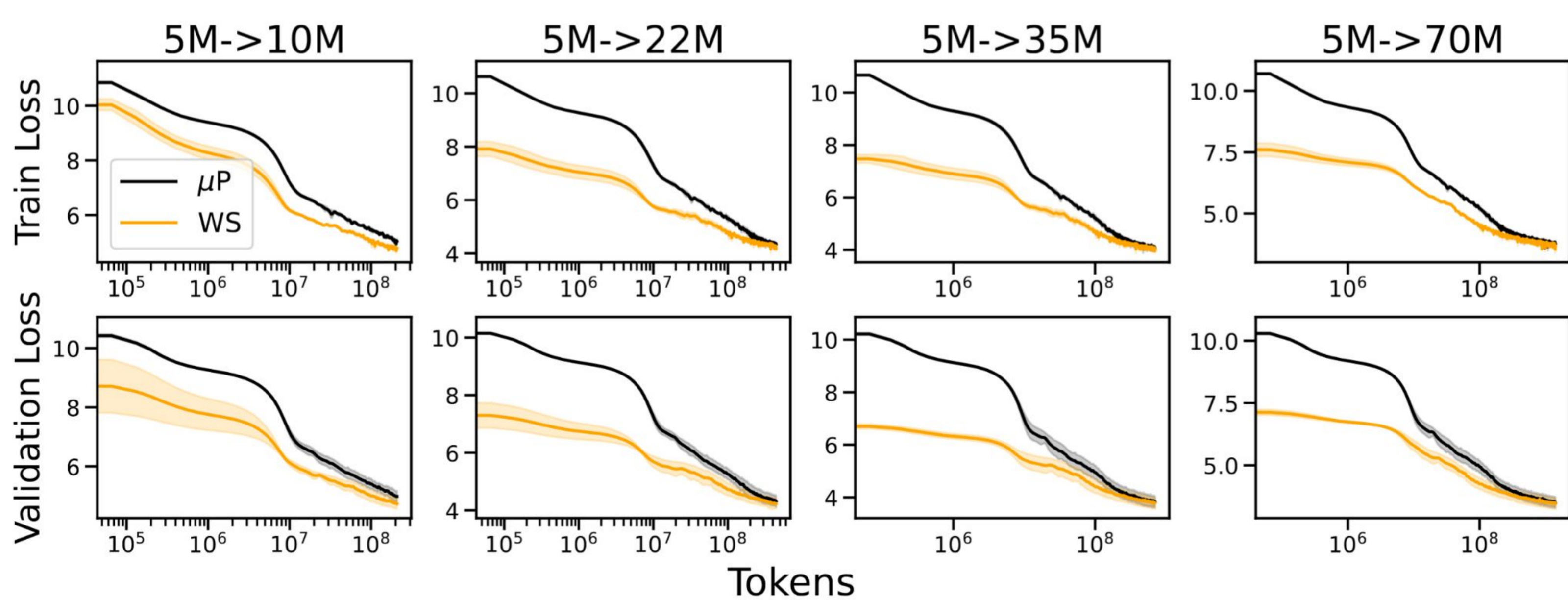
$$\theta_{target}^l = \lambda_{shrink} \cdot \text{Pad}_0(\theta_{base}^l, p, q) + \mathcal{N}(0, \sigma_{\mu P}^l)^2$$

$\lambda_{shrink} \in \mathbb{R}^1$ ; base  $\rightarrow$  smaller model; target  $\rightarrow$  larger model  
 $\theta^l \rightarrow$  weights as tensor in layer l  
 $\text{Pad}_0 \rightarrow$  zero pads a tensor to a larger tensor



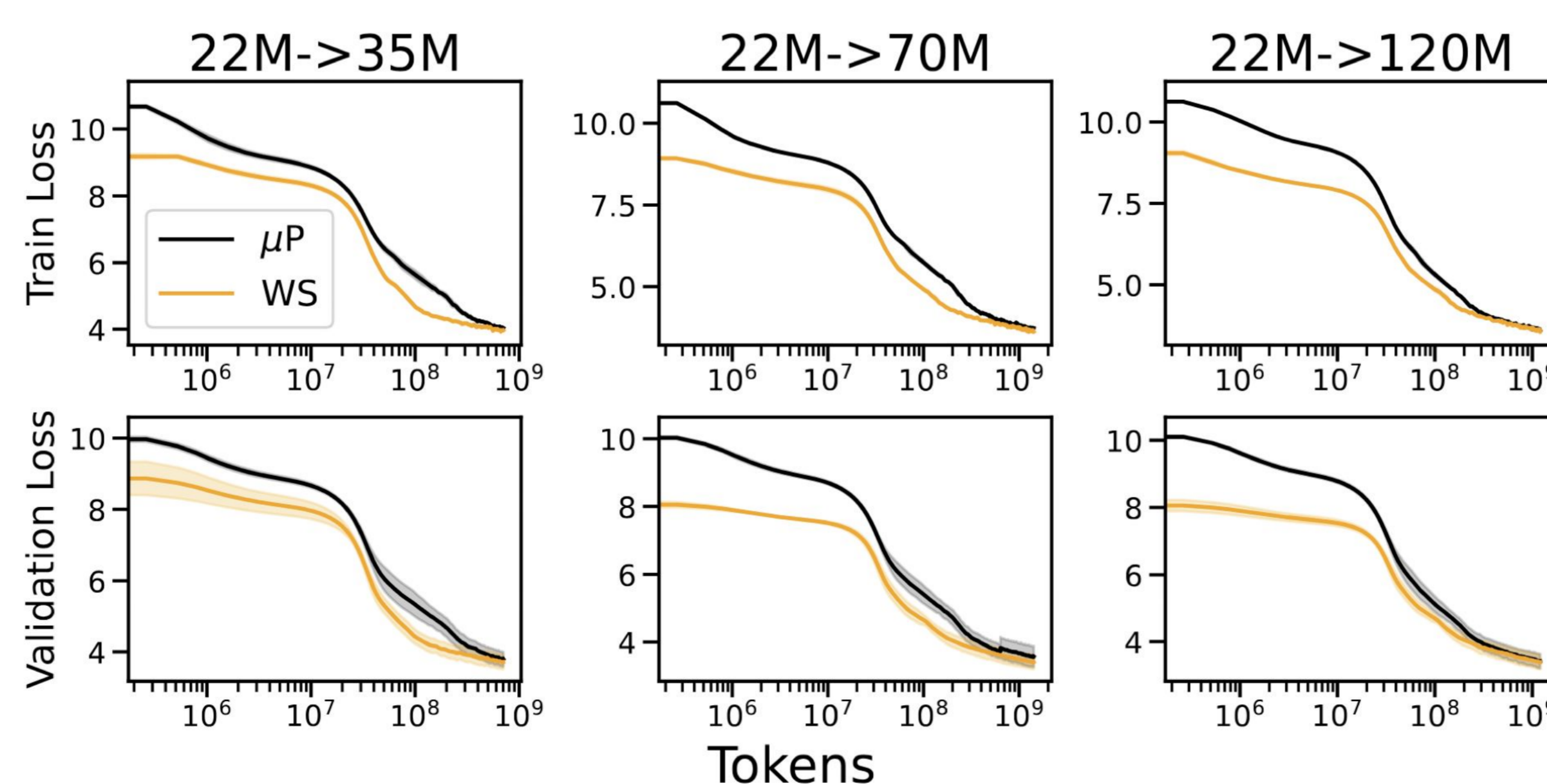
$\lambda_{shrink} = 0$  recovers  $\mu P$  without warmstarting

## Empirical Evaluation

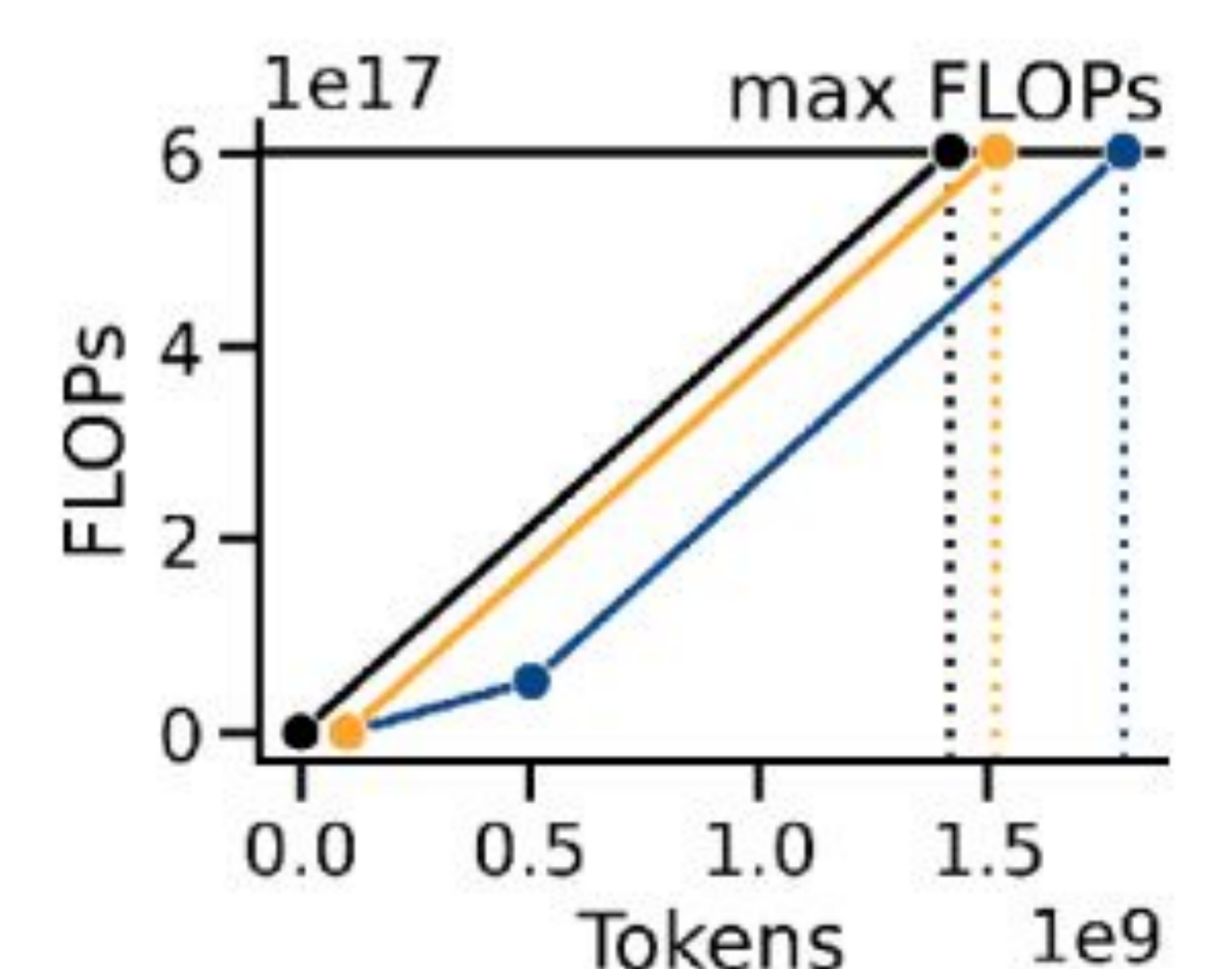


### Setup

- GPT2 architecture and vocab
- Sub-epoch training on the SlimPajama dataset [2]
- All models trained for a compute budget of 20 tokens per parameter [3]
- Trained with Adam with constant learning rate schedule [4]



### Future directions:



<sup>1</sup> G. Yang et al. Tuning large neural networks via zero-shot hyperparameter transfer. NeurIPS, 2021.  
<sup>2</sup> D. Soboleva, et al. SlimPajama: A 627B token cleaned and deduplicated version of RedPajama, HF, 2023.  
<sup>3</sup> G. J. Hoffmann et al. Training Compute-Optimal Large Language Models. 2022.  
<sup>4</sup> A Hägele et al. Scaling Laws and Compute-Optimal Training Beyond Fixed Training Durations. NeurIPS 2024.  
<sup>5</sup> G. Yang and E. J Hu. Tensor programs IV: Feature learning in infinite-width neural networks. ICML, 2021.  
<sup>6</sup> J. Ash and R. P. Adams. On warm-starting neural network training. NeurIPS, 2020.  
<sup>7</sup> K. E. Everett et al. Scaling exponents across parameterizations and optimizers. ICML, 2024.

