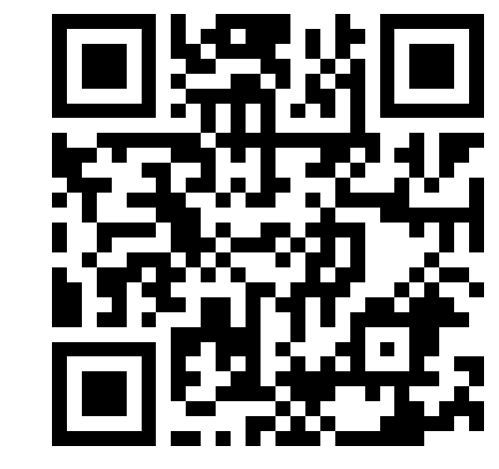


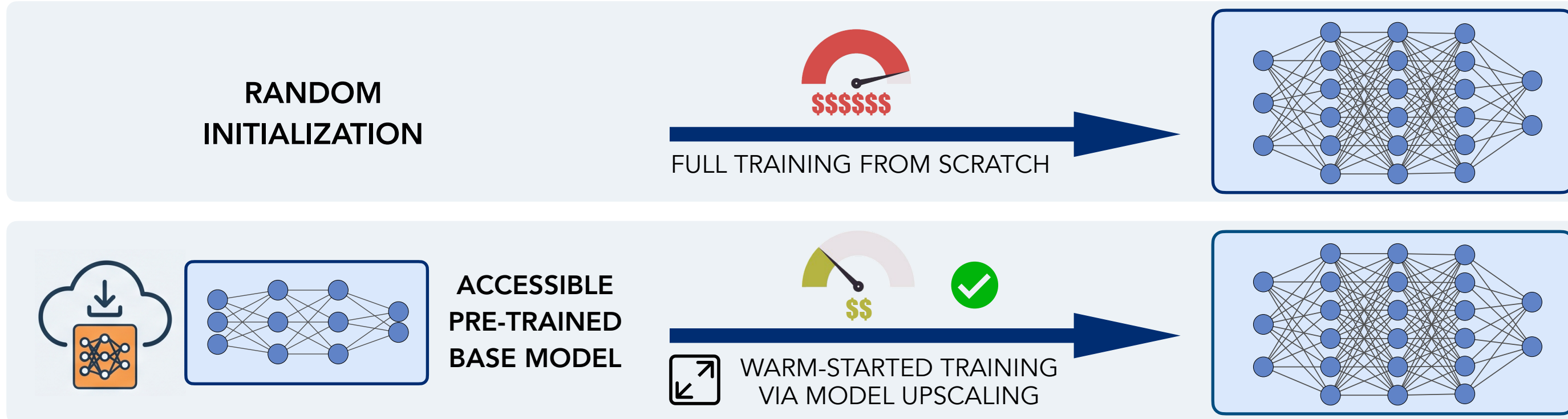


ICML



arXiv

### Model Upscaling



#### Prior work & limitations

- Based on function-preserving weight expansion [1].
- The critical gap:** Performance is sensitive to hyperparameters, and tuning at the target scale is expensive.

#### Our contribution:

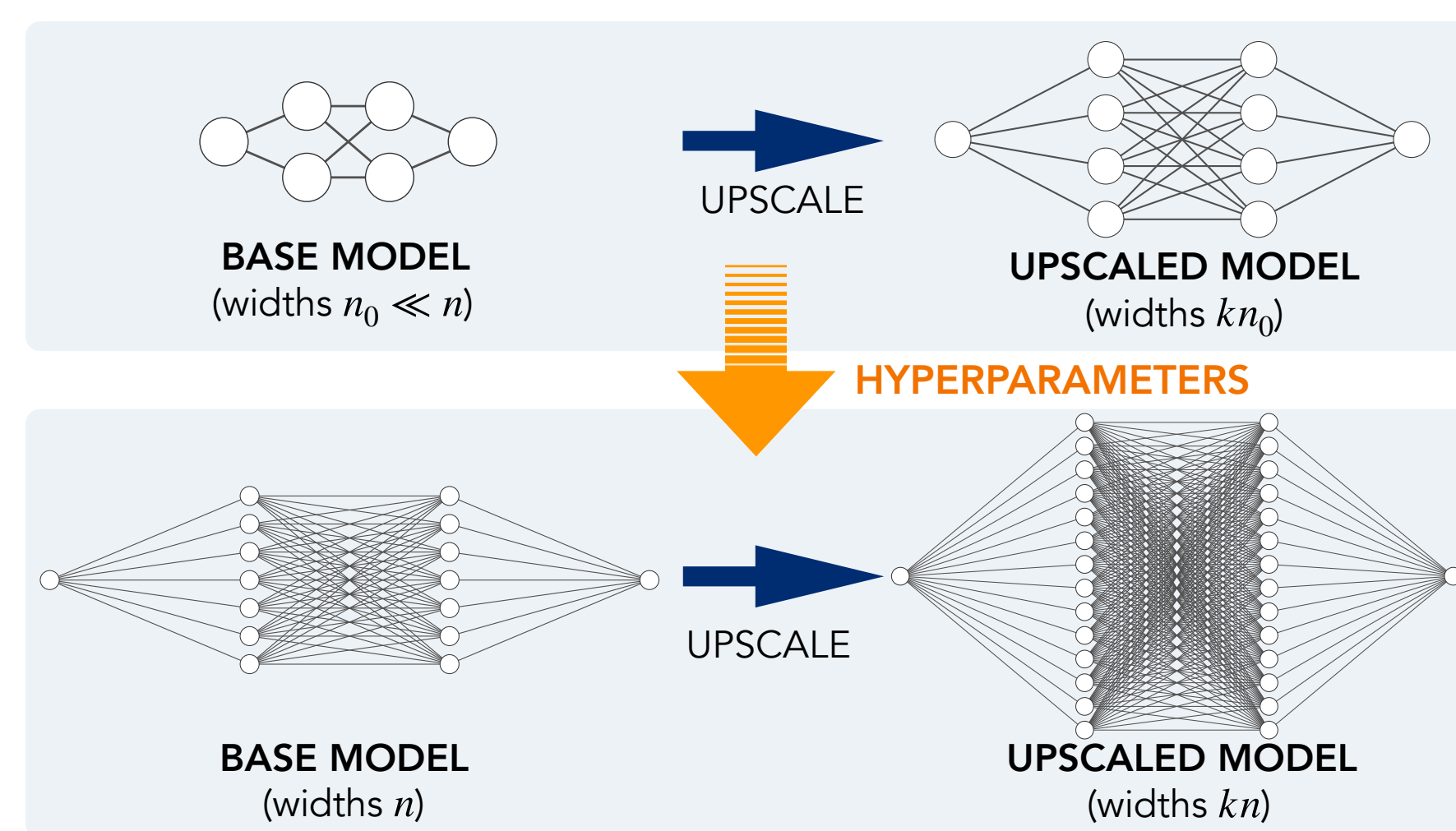
a **theoretically-grounded** upscaling method (with respect to widths) that allows **hyperparameter transfer** via  $\mu$ P.

### Proposed Upscaling Algorithm

#### Meta-algorithm 1 Upscaling procedure

**Input:** Base model checkpoint at width  $n$  with the optimizer checkpoint; expansion factor  $k$ .  
**Output:** Trained upscaled model at width  $N = nk$ .

- (HP tuning)** Tune  $\bar{\sigma}_\Delta$  and  $\bar{\gamma}^\uparrow$  by running Steps 1-4 on a small proxy ( $n_0 \rightarrow kn_0, n_0 \ll n$ ) and selecting values that minimize terminal training loss.
- Construct an equivalent wider model of width  $N$  by duplicating and rescaling the weights.
- Add noise to the weights following the  $\mu$ P initialization scaling, controlled by  $\bar{\sigma}_\Delta$ .
- Transfer the base optimizer's internal states to the upscaled model's optimizer, analogously to the weights.
- Train the upscaled model under  $\mu$ P with LR base constant  $\bar{\gamma}^\uparrow$ .

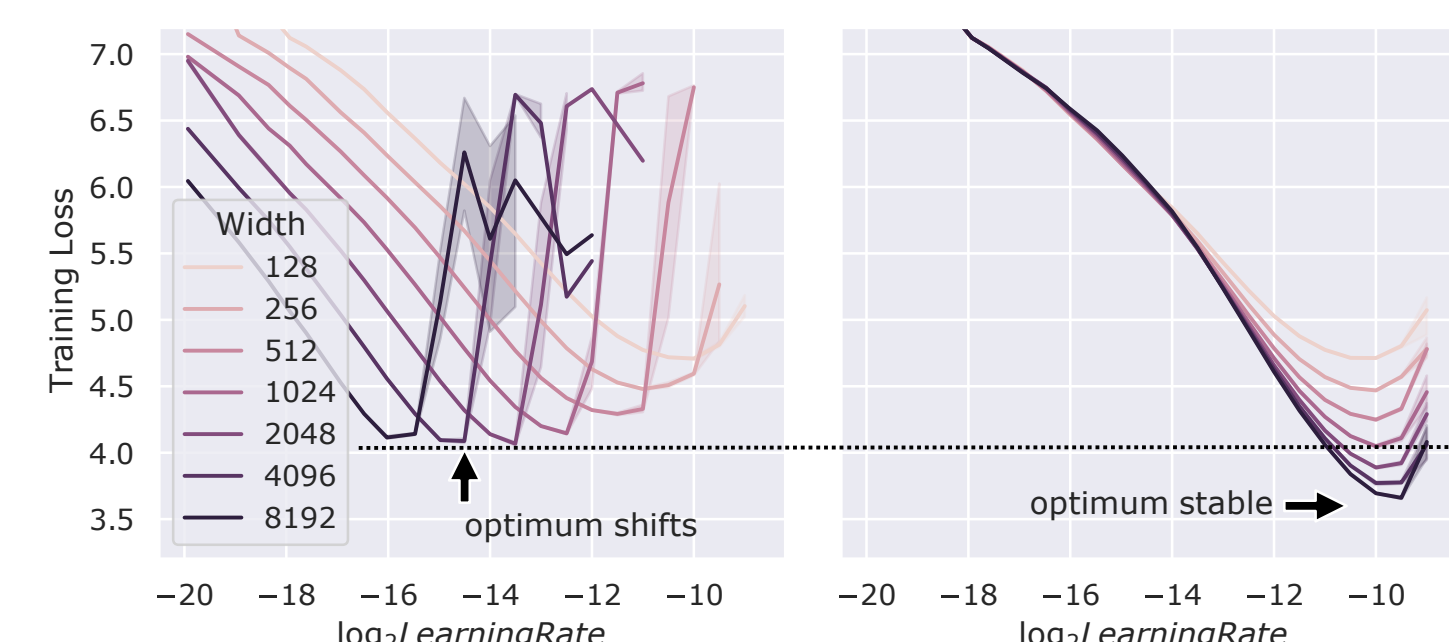


#### Maximal Update Parametrization ( $\mu$ P)

Weight type	Init Var. $\mathcal{N}(0, B \cdot \bar{\sigma}^2)$	LR $\gamma = C \cdot \bar{\gamma}$	WD (vanilla) $\lambda = D \cdot \bar{\lambda}$	WD (decoupled) $\lambda = \bar{D} \cdot \bar{\lambda}$	Extra HPs $\epsilon = E \cdot \bar{\epsilon}$	Output multiplier
Vector-like	1	$n^m$	$n^{-1}$	$n^{-m}$	$n^{-1}$	$n^{-1}$
Matrix-like	$n_{in}^{-1}$	$n_{out} n_{in}^{-1}$	$n_{in} n_{out}^{-1}$	$n_{in} n_{out}^{-m}$	$n_{out}^{-1}$	--

- $m$  depends on the optimizer:  $m = 1$  for SGD and  $m = 0$  for Adam.
- $\epsilon$  stands for extra hyperparameters that require rescaling, such as the numerical stability constant  $\epsilon_{ps}$  in Adam.

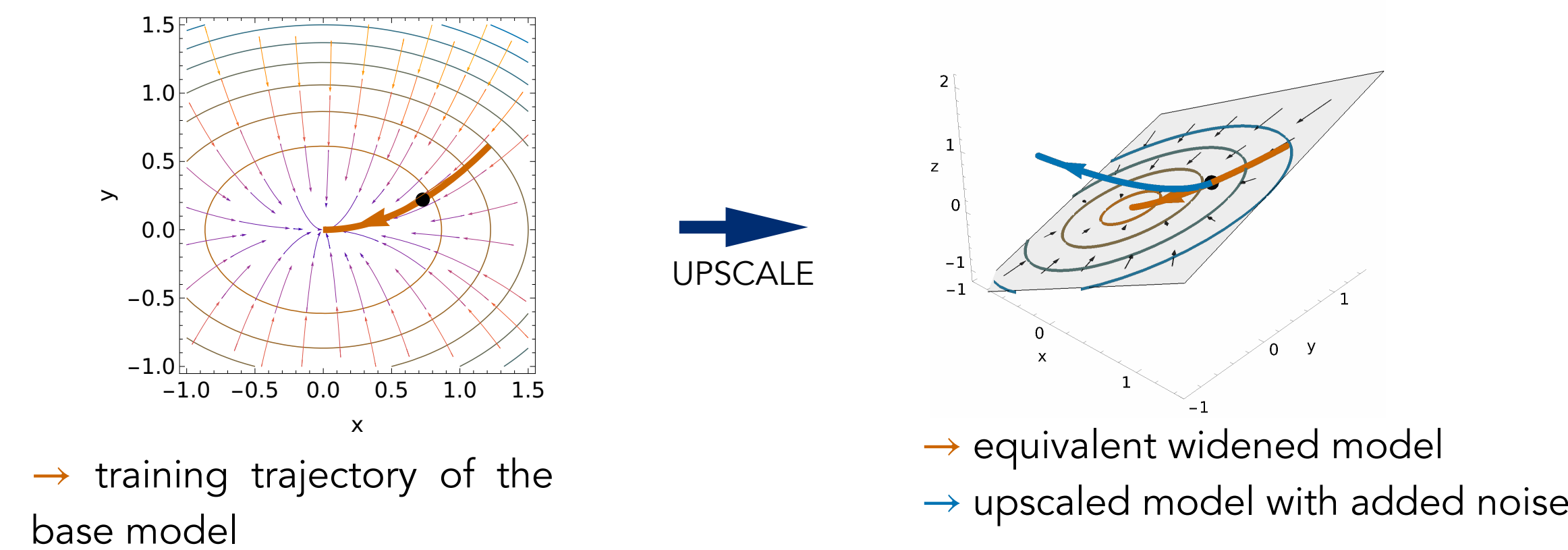
- Yields optimal feature learning behavior in the infinite-width limit [2].
- Allows zero-shot hyperparameter transfer for training from scratch [3].



### From Static Equivalence to Dynamic Equivalence

#### Theoretical guarantee 1:

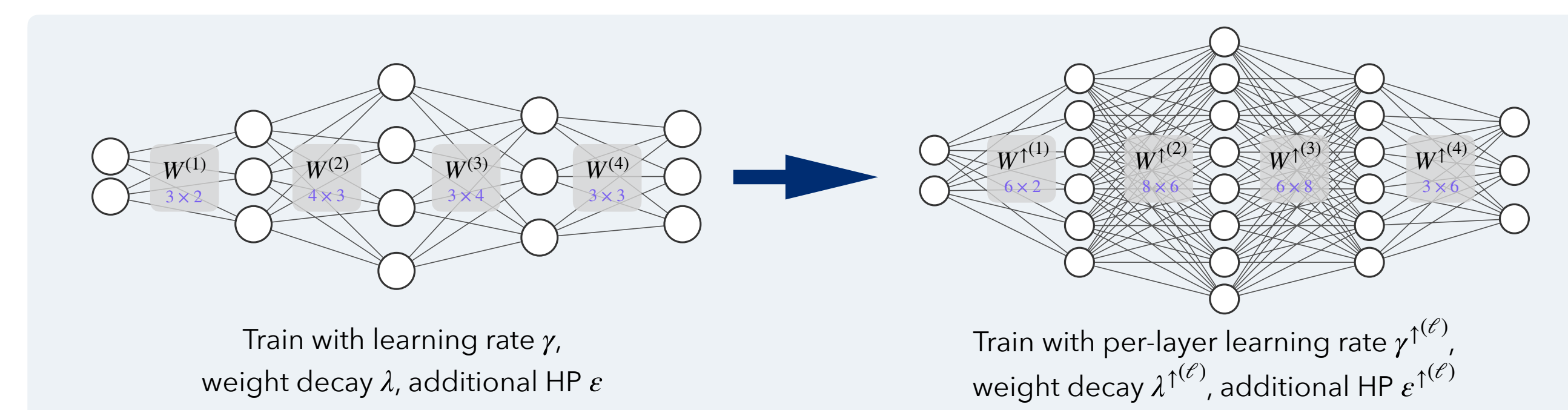
Injecting zero noise ( $\bar{\sigma}_\Delta = 0$ ) yields a widened model whose training trajectory **exactly** matches that of the base model.



**Theorem:** For standard architectures (e.g., MLP, CNN, RNN, Transformer) and optimizers (e.g., SGD, Adam), the following models are **dynamically equivalent** (expresses the same function at each training step).

- The narrow **base model** trained with given HPs.
- The **widened model** (weights constructed in the following way) trained under  $\mu$ P with the same HP base constant as the base model.

Vector-like	$n \times d \rightarrow kn \times d$	<code>W.repeat_interleave(k, dim=0)</code>
Matrix-like	$n_{out} \times n_{in} \rightarrow k_{out} n_{out} \times k_{in} n_{in}$	<code>W.repeat_interleave(k_out, dim=0).repeat_interleave(k_in, dim=1) / k_in</code>



### Infinite-Width Limit and Hyperparameter Transfer

- $\mu$ P-based HP transfer [3] applies to training from scratch only; the Tensor Programs framework [2,3] does not support the widening operation.
- We extend Tensor Programs to support mid-training upscaling.
- Key idea:** Partition the width- $kn$  activations into  $k$  groups of  $n$  neurons each. Each group uses only standard operations, placing upscaled training within the existing framework. In the infinite-width limit, activations now behave like i.i.d.  $k$ -dim Gaussian blocks with non-trivial covariance, rather than i.i.d. scalars.
- Result:** Our algorithm's **design choices** are precisely those that keep the entire training process (including mid-training upscaling) under  $\mu$ P.

#### Theoretical guarantee 2:

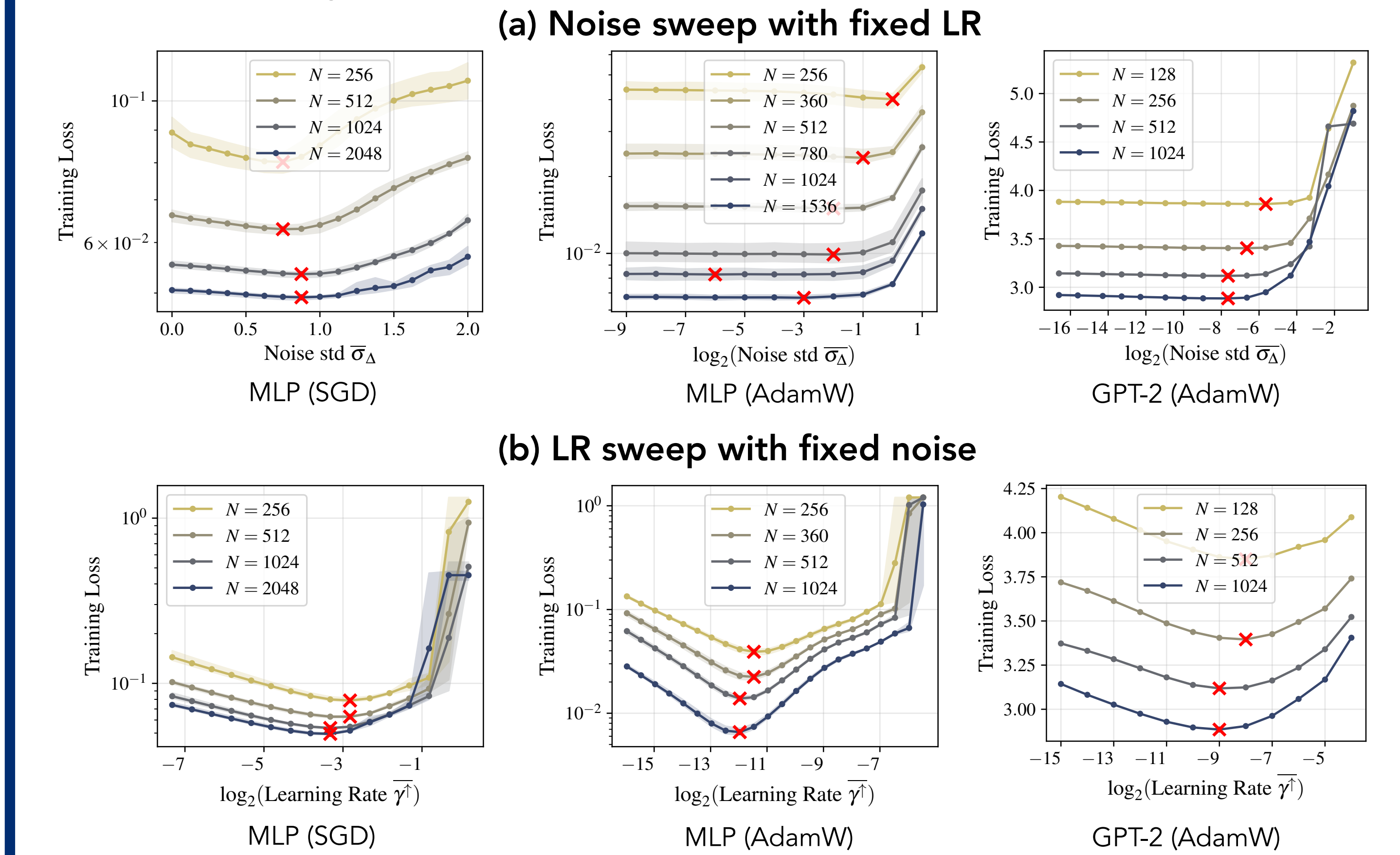
In the infinite-width limit, the injected noise is both **safe** (base-model signal preserved) and **useful** (noise's contribution neither vanishes nor explodes).

#### Theoretical guarantee 3:

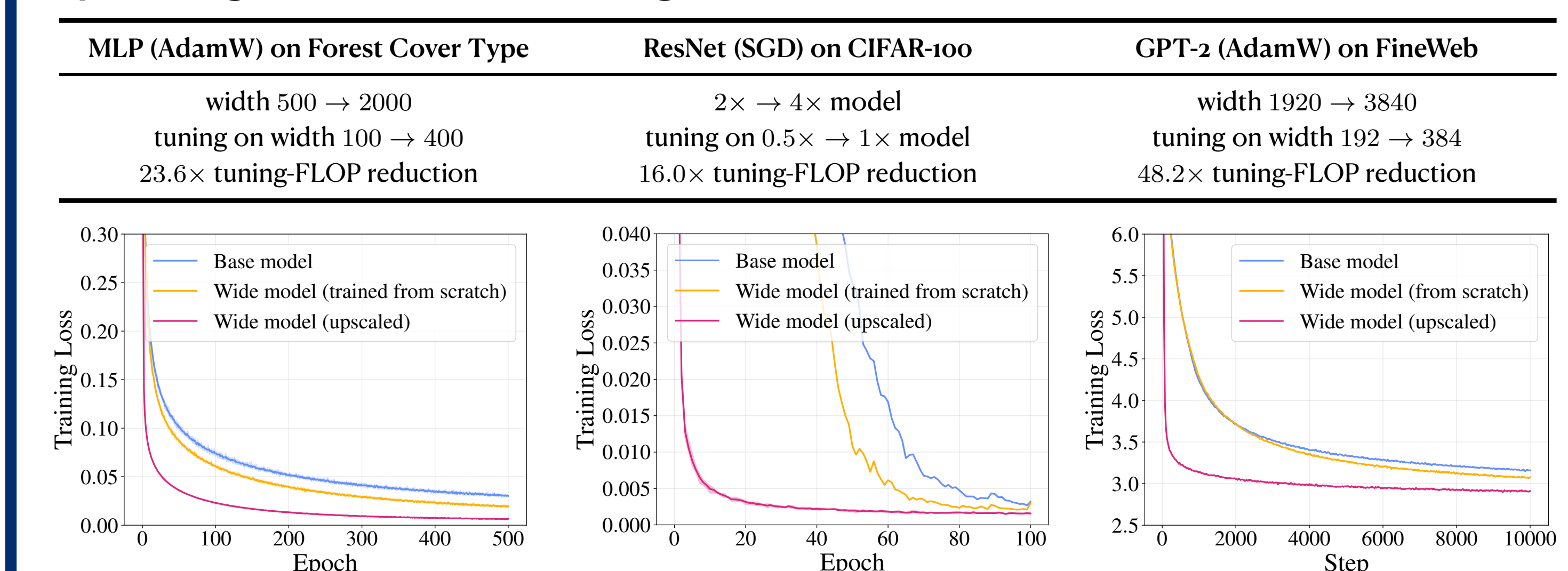
Hyperparameters  $\bar{\sigma}_\Delta$  (noise level) and  $\bar{\gamma}^\uparrow$  (upscaled model LR) can be tuned efficiently via  $\mu$ P-grounded **hyperparameter transfer**.

### Experiments

#### Validation of hyperparameter transfer



#### Upscaling accelerates training



### Conclusions and Future Work

#### Key contributions

- Provides the first rigorous HP transfer framework for model upscaling, applicable to general architectures and optimizers.
- Extends the theoretical foundation from static to dynamic equivalence.
- Connects model upscaling to the  $\mu$ P infinite-width theory.

#### Future work

- Extend to depth-wise and joint width-and-depth upscaling.
- Understand generalization behavior under upscaling (current framework addresses training dynamics only).
- Validate at industry-scale settings, where larger target-to-tuning ratio  $n/n_0$  yields greater HP tuning cost savings.

### References

- Chen, Tianqi, Ian Goodfellow, and Jonathon Shlens. "Net2Net: Accelerating Learning via Knowledge Transfer." *ICLR* (2016).
- Yang, Greg, and Edward J. Hu. "Tensor programs IV: Feature learning in infinite-width neural networks." *ICML* (2021).
- Yang, Greg, et al. "Tensor programs V: Tuning large neural networks via zero-shot hyperparameter transfer." *NeurIPS* (2021).