

Lecture 6

- RMS \rightarrow RMS norm
- Maximal update parameterization. (μP)

① Recall, induced matrix norm:

$$\|A\|_{\alpha \rightarrow \beta} = \max_{\|\vec{x}\|_\alpha = 1} \|A\vec{x}\|_\beta$$

RMS \rightarrow RMS norm \rightarrow Motivated by Xavier initialization.

$$\|A\|_{\text{RMS} \rightarrow \text{RMS}} = \max_{\|\vec{x}\|_{\text{RMS}} = 1} \|A\vec{x}\|_{\text{RMS}}$$

$$A \in \mathbb{R}^{d_{\text{out}} \times d_{\text{in}}}$$

$$d_{\text{out}} \left\{ \begin{array}{c} \boxed{\quad} \\ \vdots \end{array} \right.$$

$$A\vec{x} = h$$

$\uparrow \mathbb{R}^{d_{\text{in}}}$ $\uparrow \mathbb{R}^{d_{\text{out}}}$

$$= \max_{\|\vec{x}\|_2 = \sqrt{d_{\text{in}}}} \frac{1}{\sqrt{d_{\text{out}}}} \cdot \|A\vec{x}\|_2.$$

$$= \sqrt{\frac{d_{\text{in}}}{d_{\text{out}}}} \max_{\|\vec{x}\|_2 = 1} \|A\vec{x}\|_2$$

Spectral norm
 $\|A\|_2$!

$$= \sqrt{\frac{d_{\text{in}}}{d_{\text{out}}}} \|A\|_2$$

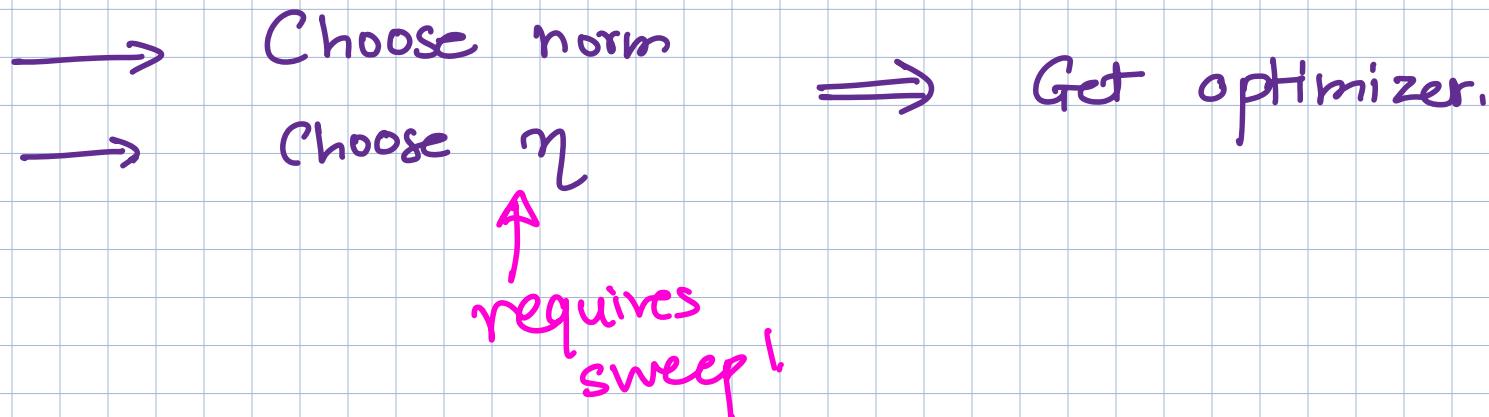
② Recall, optimizer recipe.

$$\operatorname{argmin} \|\Delta w\| \leq \eta$$

Choose appropriate norm.

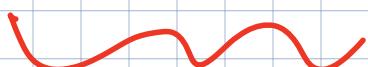
$\langle \nabla_w L(w), \Delta w \rangle$.

Gradient of Loss
change in w .



$$\operatorname{argmin} \|\Delta w\|_2 \leq \eta$$

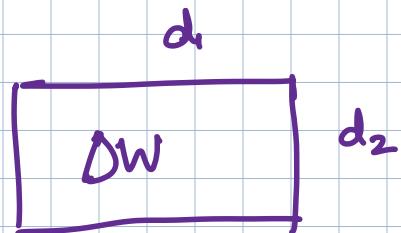
$\langle \nabla_w L(w), \Delta w \rangle$



So can we use this induced RMS \rightarrow RMS norm in our optimizer recipe?

$$\|\Delta W\|_{\text{RMS}^*} \leq \gamma ?$$

Key observation: If we choose this norm, we can focus on only one hyperparameter γ across all layers, but the effect will be layer specific learning rates, due to fan-in and fan-out.



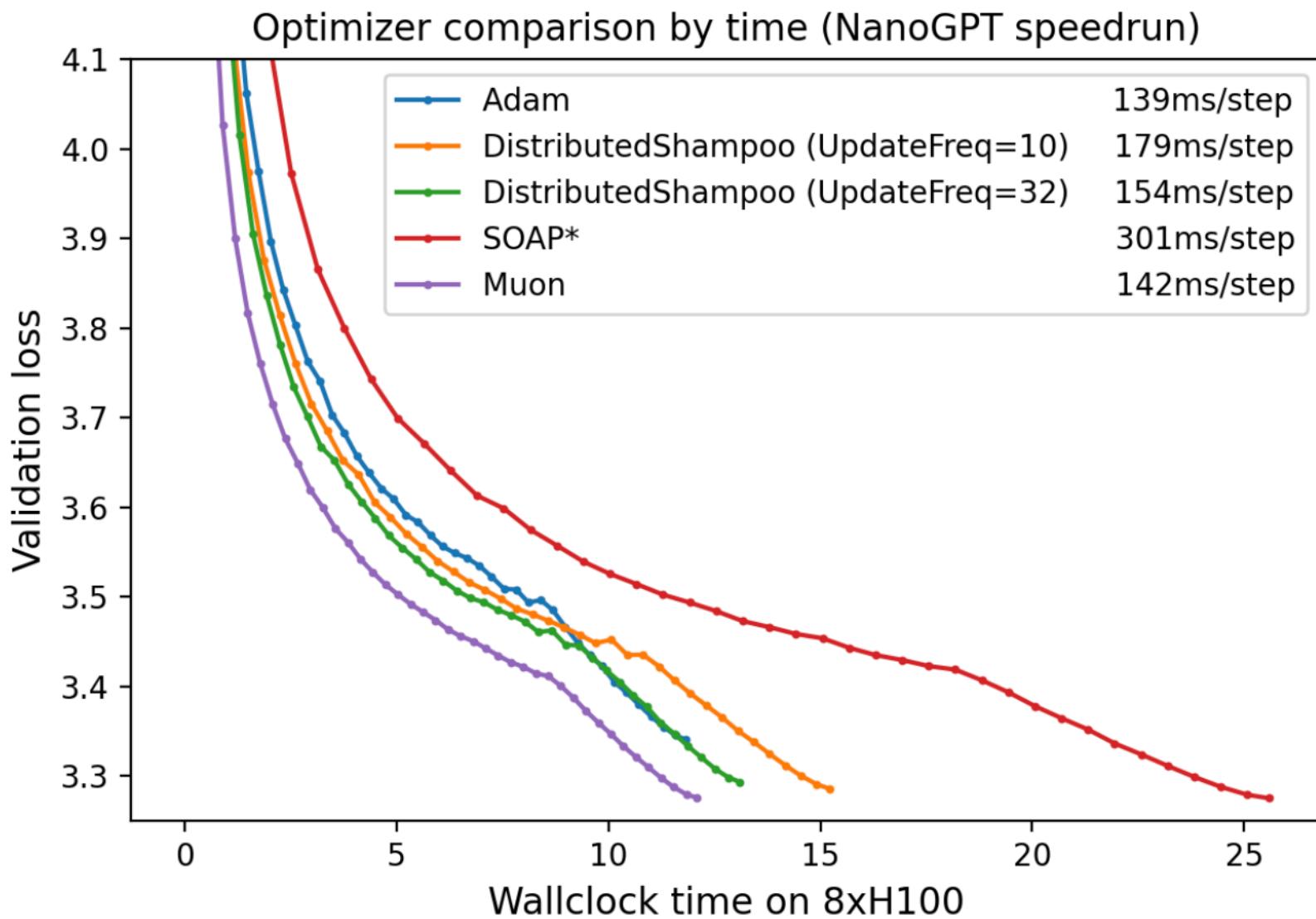
$$\begin{aligned} \|\Delta W\|_{\text{RMS}^*} \leq \gamma &\Rightarrow \sqrt{\frac{d_1}{d_2}} \|\Delta W\|_2 \leq \gamma \\ &\Rightarrow \|\Delta W\|_2 \leq \gamma \sqrt{\frac{d_1}{d_2}} \end{aligned}$$

Observe that RMS \rightarrow RMS norm is rescaling of the spectral norm.

- Use in recipe will give a scaled version of problem from last time
- leads to the Muon optimizer.

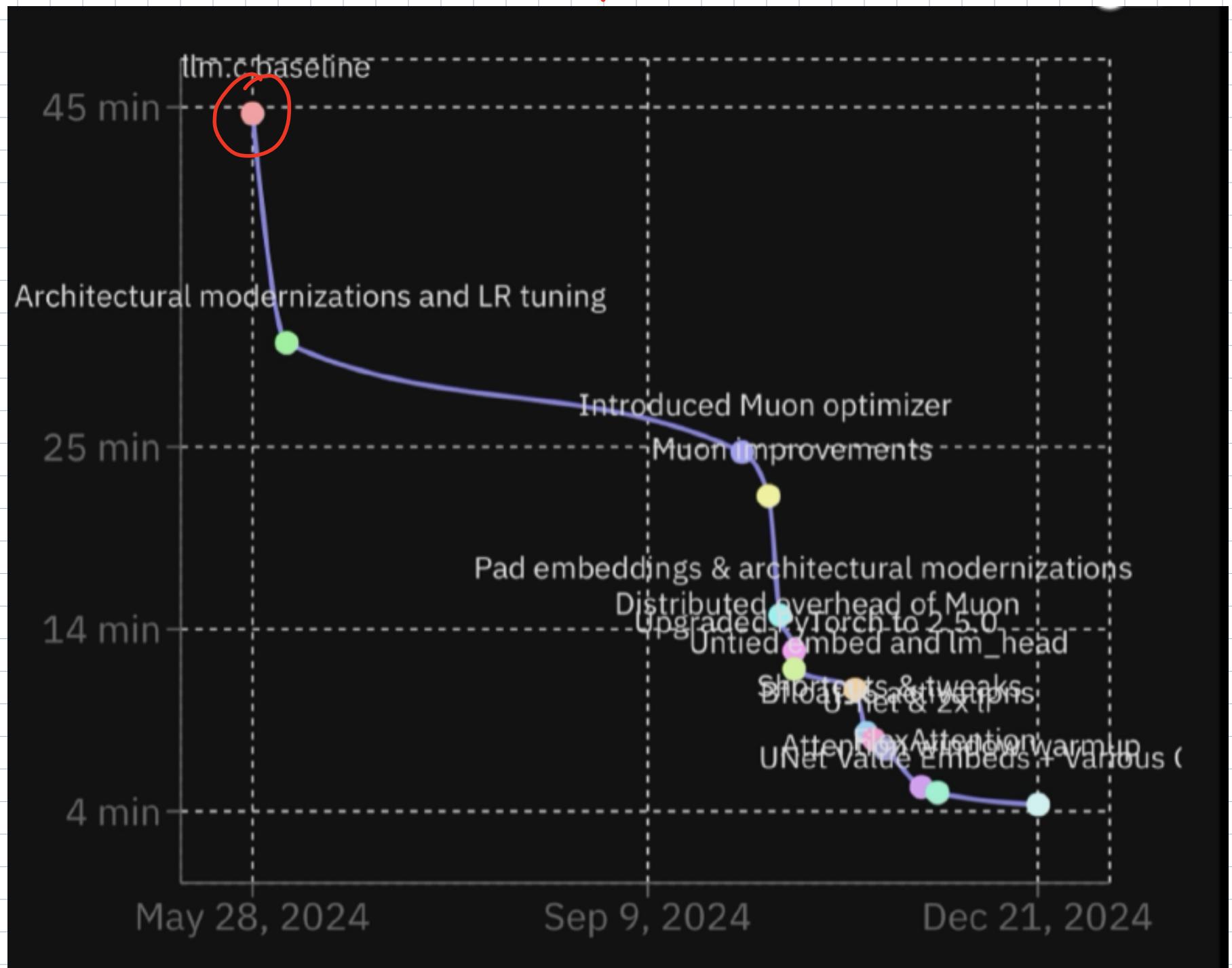
Impact.

kellerjrrdang.github.io.



*SOAP is under active development. Future versions will significantly improve the wallclock overhead.

Nano GPT speedrun



MUON IS SCALABLE FOR LLM TRAINING

TECHNICAL REPORT

Jingyuan Liu¹ Jianlin Su¹ Xingcheng Yao² Zhejun Jiang¹ Guokun Lai¹ Yulun Du¹
Yidao Qin¹ Weixin Xu¹ Enzhe Lu¹ Junjie Yan¹ Yanru Chen¹ Huabin Zheng¹
Yibo Liu¹ Shaowei Liu¹ Bohong Yin¹ Weiran He¹ Han Zhu¹ Yuzhi Wang¹
Jianzhou Wang¹ Mengnan Dong¹ Zheng Zhang¹ Yongsheng Kang¹ Hao Zhang¹
Xinran Xu¹ Yutao Zhang¹ Yuxin Wu¹ Xinyu Zhou¹* Zhilin Yang¹

¹ Moonshot AI ² UCLA



22 February 2025

Different approach.

Commit to sign SGD/Adam.

$$w_{t+1} = w_t - \eta \cdot \text{sign}(\nabla_w L(w)).$$

But I want to rethink step size choice

Factor in that different layers have different architectures.

Maybe RMS norm can help?

Want : $\|\Delta w\|_{\text{RMS} \rightarrow \text{RMS}} \leq \gamma$

Consider mini-batch size 1.

$\Rightarrow \nabla_w L$ gradient is low-rank, specifically rank 1.

$$\Rightarrow \nabla_w L = \sigma \cdot \vec{u} \vec{v}^T \quad \left| \begin{array}{l} W\vec{x} = \vec{h} \\ \frac{\partial L}{\partial w} = \frac{\partial h}{\partial w} \frac{\partial L}{\partial h} \end{array} \right.$$

What is the spectral norm? Frobenius norm?

$$\|\nabla_w L\|_2 = \sigma \quad \|\nabla_w L\|_F = \sigma$$

Consider $\text{sgn}(\nabla_w L) = S$. (elementwise)

What is rank S ?

$$\begin{aligned} S &= \text{sgn}(\sigma \vec{u} \vec{v}^T) = \text{sgn}(\vec{u} \vec{v}^T) \\ &= \text{sgn}(\vec{u}) \cdot \text{sgn}(\vec{v}^T) \end{aligned}$$

$$\begin{aligned} \text{sgn}(u_i v_j) &= \text{sgn}(u_i) \text{sgn}(v_j) \end{aligned}$$

\Rightarrow Rank S must be 1

Consider mini-batch size 1.

$\Rightarrow \nabla_w L$ gradient is low-rank, specifically rank 1.

$$\Rightarrow \nabla_w L = \sigma \cdot \vec{u} \vec{v}^T$$

$$\|\nabla_w L\|_2^2 = \sigma^2$$

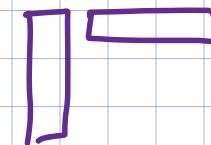
$$\|\nabla_w L\|_F^2 = \text{trace}(\sigma^2 \vec{v} \vec{u}^T \vec{u} \vec{v}^T) = \sigma^2 \text{ (sum of singular values squared)}$$

Consider $\text{sgn}(\nabla_w L) = S$. (elementwise)

What is rank S ?

$$\begin{aligned}\text{sgn}(S) &= \text{sgn}(\sigma \vec{u} \vec{v}^T) \\ &= \text{sgn}(\vec{u} \vec{v}^T) \\ &= \text{sgn}(\vec{u}) \cdot \text{sgn}(\vec{v}^T)\end{aligned}$$

$$\text{sgn}(u_i; v_j) = \text{sgn}(u_i) \cdot \text{sgn}(v_j)$$



So S is also rank 1.

$$W_{t+1} = W_t - \eta \cdot \underbrace{S}_{\Delta W} \cdot$$

Recall, we want DW to be small.

$$\|DW\|_{RMS \rightarrow RMS} \leq \gamma$$

$$\eta \cdot \|S\|_{RMS \rightarrow RMS} \leq \gamma \Rightarrow \eta \cdot \sqrt{\frac{d_{in}}{d_{out}}} \cdot \|S\|_2 \leq \gamma$$

What is $\|S\|_2$?

Example:

$$A = \begin{bmatrix} +2 & -1 \\ -3 & +1 \end{bmatrix}$$

$$P = \text{Sgn}(A) = \begin{bmatrix} 1 & -1 \\ -1 & 1 \end{bmatrix} \}_{2^2} \quad (-1)^2 = 1^2 = 1$$

$$\textcircled{B} \quad \|P\|_F^2 = \cancel{2 \times 2 \times 1} = 4$$

$$\|A\|_F^2 = \sum_{\text{Singu. values}} \sigma_i^2 = \text{trace}(A^T A) = \text{trace}(A A^T) = \sum_{i,j} a_{ij}^2$$

Recall, we want ΔW to be small.

$$\|\Delta W\|_{RMS \rightarrow RMS} \leq \gamma$$

$$\eta \cdot \|S\|_{RMS \rightarrow RMS} \leq \gamma \Rightarrow \eta \cdot \sqrt{\frac{d_{in}}{d_{out}}} \cdot \|S\|_2 \leq \gamma$$

What is $\|S\|_2$?

S is rank 1, so spectral norm = Frobenius norm.

$$\|S\|_F^2 = d_{in} \times d_{out}$$

Entries are just ± 1 .

$$\Rightarrow \|S\|_2 = \sqrt{d_{in} \times d_{out}}$$

$$\Rightarrow \eta \cdot \sqrt{\frac{d_{in}}{d_{out}}} \cdot \sqrt{d_{in} \times d_{out}} \leq \gamma$$

$$\Rightarrow \eta \leq \frac{1}{d_{in}} \cdot \gamma$$

Learning rate can scale with layer.

So my Adam / signSGD learning rate should scale as $\frac{1}{d\eta}$

to maintain $\|\Delta W\|_{\text{RMS} \rightarrow \text{RMS}} \leq \gamma$

→ This captures an essence of the μP parameterization
(Maximal update parameterization)

Yang et al. 2022.

Notes:

① As batch size increases, Frobenius norm will only be an upper bound to spectral norm.

So our bound will be conservative.

GPT-3, 2020.

Model Name	n_{params}	n_{layers}	d_{model}	n_{heads}	d_{head}	Batch Size	Learning Rate
GPT-3 Small	125M	12	768	12	64	0.5M	6.0×10^{-4}
GPT-3 Medium	350M	24	1024	16	64	0.5M	3.0×10^{-4}
GPT-3 Large	760M	24	1536	16	96	0.5M	2.5×10^{-4}
GPT-3 XL	1.3B	24	2048	24	128	1M	2.0×10^{-4}
GPT-3 2.7B	2.7B	32	2560	32	80	1M	1.6×10^{-4}
GPT-3 6.7B	6.7B	32	4096	32	128	2M	1.2×10^{-4}
GPT-3 13B	13.0B	40	5140	40	128	2M	1.0×10^{-4}
GPT-3 175B or "GPT-3"	175.0B	96	12288	96	128	3.2M	0.6×10^{-4}

Table 2.1: Sizes, architectures, and learning hyper-parameters (batch size in tokens and learning rate) of the models which we trained. All models were trained for a total of 300 billion tokens.

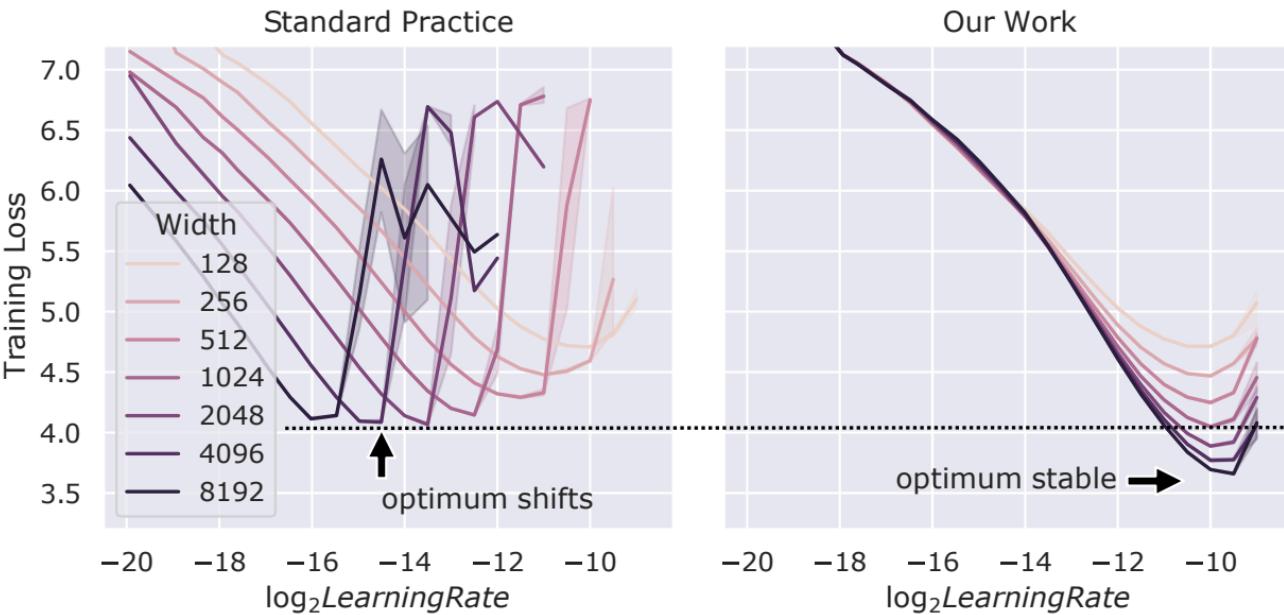


Figure 1: Training loss against learning rate on Transformers of varying d_{model} trained with Adam. Conventionally and in contrast with our technique, different widths do not share the same optimal hyperparameter; wider networks do not always perform better than narrower ones; in fact they underperform the same-width networks in our technique even after tuning learning rate (see dashed line). See Sections 3 and 4 for experimental setup.

Yang et al 2022.

What is parameterization*? Why does it matter?

- About choosing the right "units"
- A key challenge for neural net training is hyperparameter search.
- Particularly bad if the network is large.
 - Experimenting is very expensive.
- Can we optimize hyperparameters on smaller networks and then "transfer" them to larger ones? What is the right scaling?

* UP:
How to change
hyp. params when NN
width changes.

Example: x_1, x_2, \dots, x_n iid from zero-mean, unit variance distribution.

Then $\frac{1}{\sqrt{n}} (x_1 + x_2 + \dots + x_n) \xrightarrow[n \rightarrow \infty]{} N(0, 1)$

But $\frac{1}{n} (x_1 + x_2 + \dots + x_n) \xrightarrow[n \rightarrow \infty]{} 0$

$$\frac{1}{1} (x_1 + x_2 + \dots + x_n) \xrightarrow[n \rightarrow \infty]{} \infty$$

$\frac{1}{\sqrt{n}}$ is the right order of "Scaling factor" so that $c_n(\sum x_i)$ converges to something non-trivial.

Consider

$$F_n(c) \doteq \mathbb{E}_{x_1, x_2, \dots, x_n} f(c(x_1 + x_2 + \dots + x_n))$$

$c \in \mathbb{R}$, $f: \mathbb{R} \rightarrow \mathbb{R}$ bounded, continuous.

Reparameterize $c = \alpha/\sqrt{n}$.

only function of α .

$$G_n(\alpha) \doteq F_n\left(\frac{\alpha}{\sqrt{n}}\right) \rightarrow \mathbb{E}(f(N(0, \alpha^2)))$$

as $n \rightarrow \infty$

$\alpha_n^* = \underset{\alpha}{\operatorname{argmin}} G_n(\alpha) \rightarrow$ Then $\alpha_n^* \approx \alpha_N^*$ for large enough n, N

→ We can copy α_n^* from a smaller n to a larger N .

If instead $C_n^* = \underset{C}{\operatorname{argmin}} F_n(C) \rightarrow$ cannot transfer parameters in the same way.

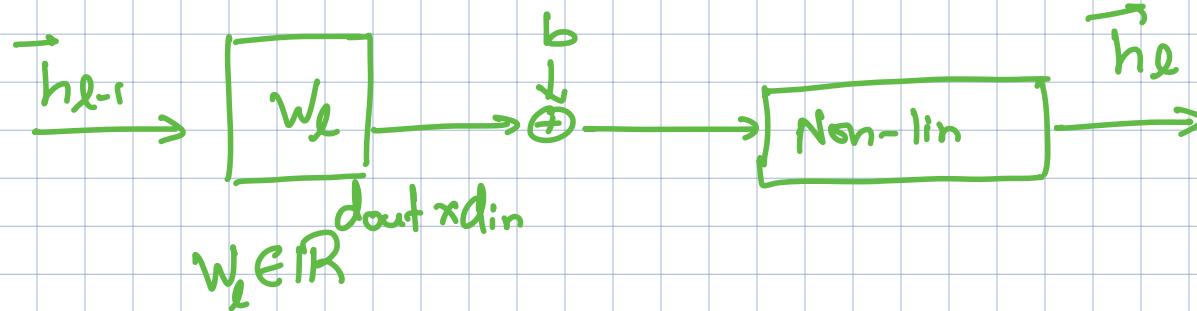
Can we do something similar for NN hyperparameters?

$$\eta \leq \frac{1}{d_{in}} \cdot ? \quad ? \text{ still a hyperparameter.}$$

the "right scaling".

Conditions for Feature Learning

\vec{h}_l : hidden layer l .



Focus on scaling with network width.

Xavier initialization thinking suggests we want

$$\|\vec{h}_l\|_{\text{RMS}} = \frac{1}{\sqrt{d_{\text{out}}}} \|\vec{h}_l\|_2$$

$$\Rightarrow \|\vec{h}_l\|_2 \approx \sqrt{d_{\text{out}}}$$

$$\|\vec{h}_l\|_{\text{RMS}} \approx 1.$$

Relaxation of all entries approx 1.

$$\|\vec{h}_l\|_{\text{RMS}} = \Theta(1)$$

Also, want to ensure that our updates to the layers are not negligible as width grows.

$$\|\Delta \vec{h}_l\|_{\text{RMS}} = \Theta(1)$$

Desiderata: (1) $\|\vec{h}_l\|_{\text{RMS}} = \Theta(1)$

(2) $\|\Delta \vec{h}_l\|_{\text{RMS}} = \Theta(1)$

How can we achieve this?

For now, ignore bias + non-linearity.

$$\vec{h}_l = w_l \vec{h}_{l-1}$$

\Rightarrow if \vec{h}_{l-1} aligns with top singular vector of w_l . [Can show this is the case]

$$\|\vec{h}_l\|_2 = \|w_l\|_2 \cdot \|\vec{h}_{l-1}\|_2$$

$$\Theta(\sqrt{d_{out}}) = \|w_l\|_2 \Theta(\sqrt{d_{in}})$$

$$\Rightarrow \|w_l\|_2 = \Theta\left(\sqrt{\frac{d_{out}}{d_{in}}}\right) \text{ i.e. } \|w_l\|_{RMS \rightarrow RMS} = \Theta(1)$$

i.e. fixed RMS norm as

we started with!

Thus, if $\|w_l\|_{RMS \rightarrow RMS} = \Theta(1)$ and $\|\vec{h}_l\|_{RMS} = \Theta(1)$, then

$$\|\vec{h}_l\|_{RMS} = \Theta(1)$$

What about $\|\Delta h_l\|_{RMS}$?

\rightarrow Need one more condition.

$$\|\Delta \mathbf{W}_e\|_{RMS} = \Theta(1) \Rightarrow \|\Delta \mathbf{W}_e\|_2 = \Theta\left(\sqrt{\frac{d_{out}}{d_{in}}}\right)$$

$$\overrightarrow{h_e} = \mathbf{w}_e \overrightarrow{h_{e-1}}$$

$$\Rightarrow \overrightarrow{\Delta h_e} = \Delta \mathbf{w}_e \cdot \overrightarrow{h_{e-1}} + \mathbf{w}_e \cdot \overrightarrow{\Delta h_{e-1}}$$

$$\Rightarrow \|\Delta \overrightarrow{h_e}\|_2 \leq \|\Delta \mathbf{w}_e\|_2 \cdot \|\overrightarrow{h_{e-1}}\|_2 + \|\mathbf{w}_e\|_2 \|\Delta \overrightarrow{h_{e-1}}\|_2.$$

$$= \boxed{?} \quad \Theta(\sqrt{d_{in}}) + \Theta\left(\sqrt{\frac{d_{out}}{d_{in}}}\right) \cdot \Theta(\sqrt{d_{in}})$$

$$\Rightarrow \|\Delta \mathbf{h}_e\|_{RMS} = \Theta(1)$$

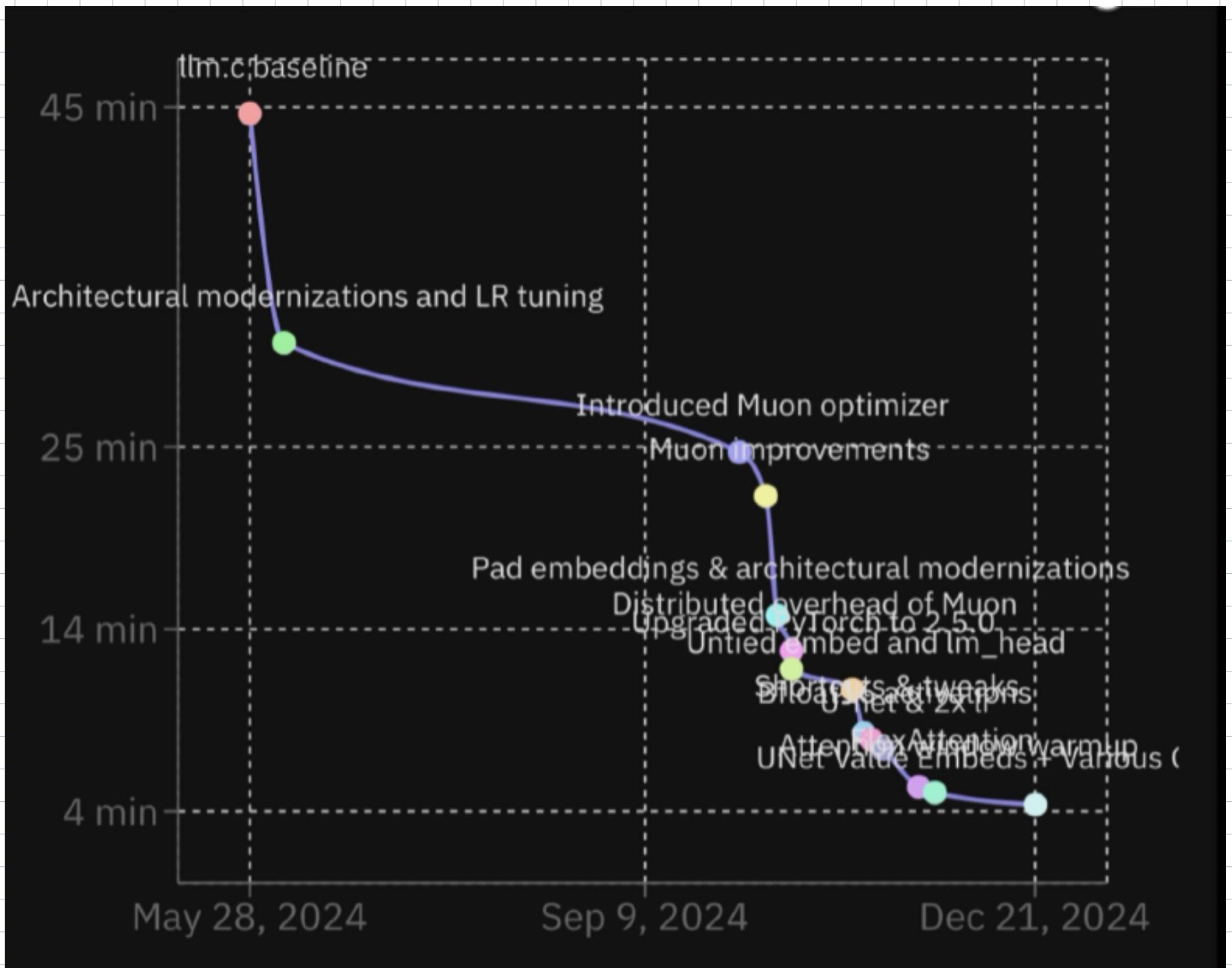
Conditions: (1) $\|\mathbf{W}\|_{RMS} = \Theta(1)$

(2) $\|\Delta \mathbf{W}\|_{RMS} = \Theta(1)$

$$\vec{\Delta h}_e = D\vec{w}_e \cdot \vec{h}_{e-1} + \vec{w}_e \cdot \vec{\Delta h}_{e-1}$$

$$\begin{aligned}\|\vec{\Delta h}_e\|_2 &= \|\vec{\Delta w}_e \vec{h}_{e-1} + \vec{w}_e \vec{\Delta h}_{e-1}\|_2 \\ &\leq \|D\vec{w}_e\|_2 \|\vec{h}_{e-1}\|_2 + \|\vec{w}_e\|_2 \|\vec{\Delta h}_{e-1}\|_2\end{aligned}$$

→ apply triangle inequality
+ sub-multiplicativity property



Consider mini-batch signSGD.

Batch size 1..

$\nabla_w L(w)$ must be in the span of the datapoint.

het $\vec{w} \vec{x} = \vec{h}$

$$w = \begin{bmatrix} w_{ij} \end{bmatrix} \quad \begin{bmatrix} w_{11} & w_{12} \\ w_{21} \end{bmatrix}$$

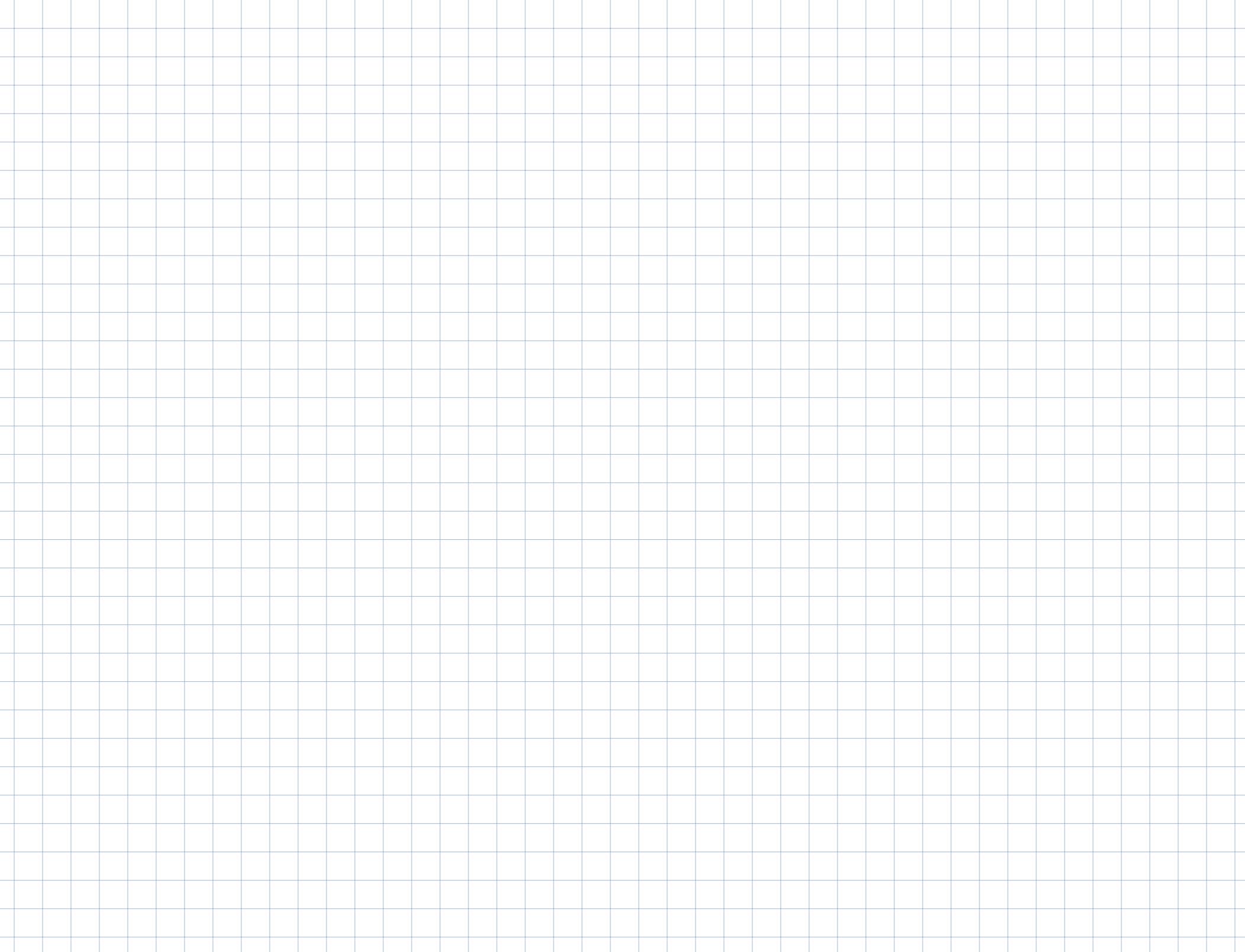
$$\frac{\partial L}{\partial w} = \frac{\partial h}{\partial w} \cdot \frac{\partial L}{\partial h}$$

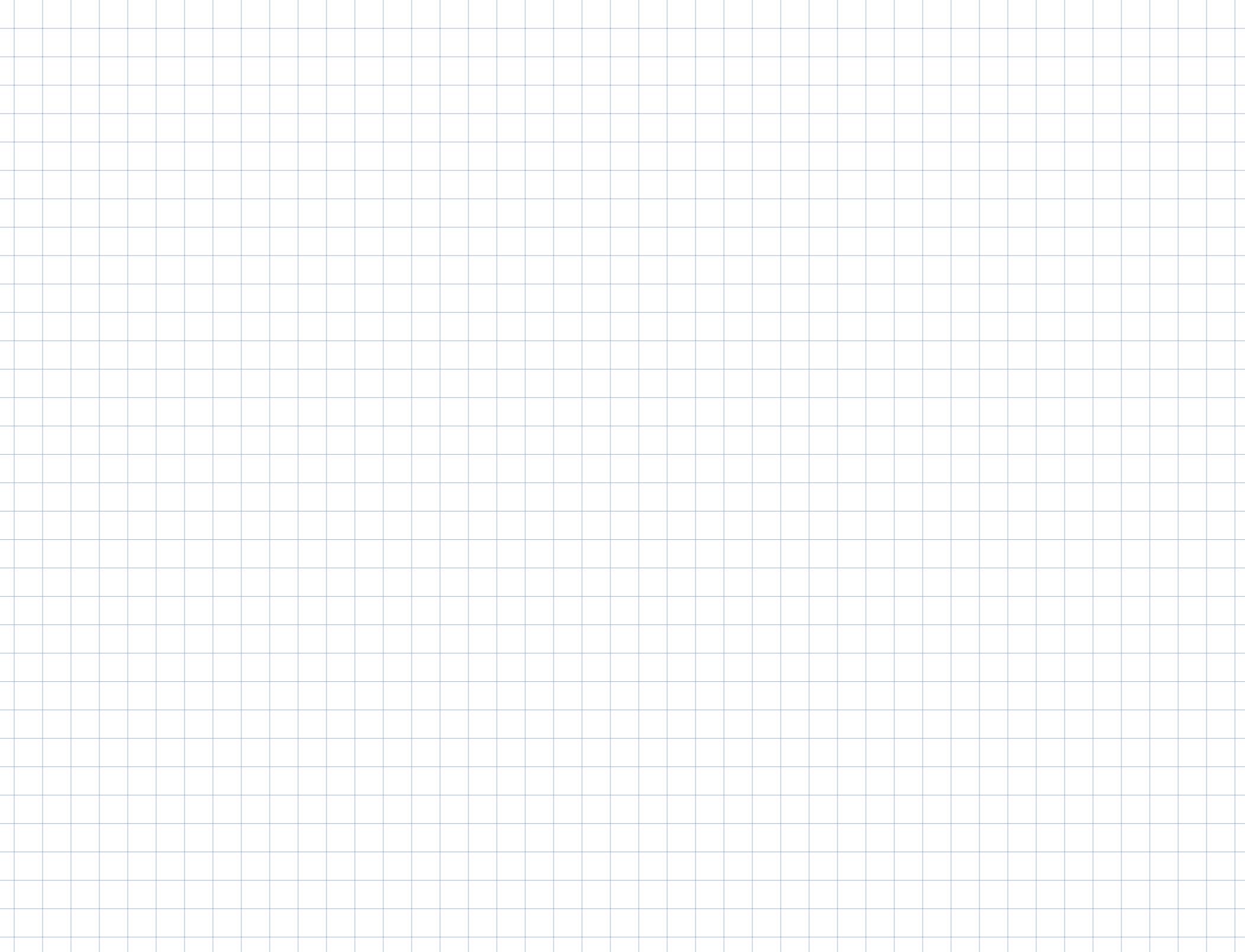
$$h_i = \sum_j w_{ij} \cdot x_j$$

A diagram illustrating the components of the gradient. On the left, a vertical vector labeled \vec{x} is shown. To its right, a horizontal arrow labeled $\frac{\partial L}{\partial h}$ points to the right.

$$\frac{\partial h_i}{\partial w_{ij}} = x_j$$

$$\frac{\partial L}{\partial w} = \vec{x} \frac{\partial L}{\partial h}$$





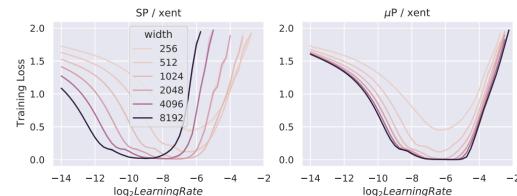


Figure 3: MLP width different hidden sizes trained for 20 epoch on CIFAR-10 using SGD. **Left** uses standard parametrization (SP); **right** uses maximal update parametrization (μP). μP networks exhibit better learning rate stability than their SP counterparts.