

# Beyond MuP: 3. Special Cases Require Special Treatment

Jianlin Su

2026-03-02

## Abstract

This article explores why certain components like Embedding layers and LM Heads require special treatment in optimization, despite being matrix-valued like other layers. Using the three stability metrics introduced in earlier articles, we analyze the initialization and steepest descent directions for these special cases and show why they differ from the Muon optimizer's approach.

## 1 Recap

In the first article [Beyond MuP: 1. Three Features of a Good Model](#), we introduced three stability metrics:

$$\text{Forward Stability: } \max_{\mathbf{x}} \|\mathbf{f}(\mathbf{x}; \boldsymbol{\omega})\|_{\text{RMS}} = \Theta(1) \quad (1)$$

$$\text{Dependency Stability: } \max_{\mathbf{x}_1, \mathbf{x}_2} \|\mathbf{f}(\mathbf{x}_1; \boldsymbol{\omega}) - \mathbf{f}(\mathbf{x}_2; \boldsymbol{\omega})\|_{\text{RMS}} = \Theta(1) \quad (2)$$

$$\text{Update Stability: } \max_{\mathbf{x}} \|\mathbf{f}(\mathbf{x}; \boldsymbol{\omega} + \Delta\boldsymbol{\omega}) - \mathbf{f}(\mathbf{x}; \boldsymbol{\omega})\|_{\text{RMS}} = \Theta(1) \quad (3)$$

These metrics share a common format: compute the RMS of the output, then take the maximum over inputs. Here,  $\mathbf{x}$  is the input,  $\boldsymbol{\omega}$  the parameters, and  $\mathbf{f}(\mathbf{x}; \boldsymbol{\omega})$  can represent a layer, block, or the entire model depending on our ability to compute the maximum.

In the previous article [Beyond MuP: 2. Linear Layers and Steepest Descent](#), we showed how to compute these metrics for linear layers by adding In Norms, and derived the Muon optimizer using the steepest descent principle.

The key contribution of the "Beyond MuP" series is not the steepest descent itself, but the identification of suitable stability metrics for arbitrary layers.

## 2 Embedding Layers

Consider an Embedding layer: input is an index  $i$ , output is the corresponding vector  $\mathbf{f}(i; \mathbf{E}) = \mathbf{E}_i$ , where  $\mathbf{E}$  is a  $|V| \times d$  matrix and  $\mathbf{E}_i \triangleq \mathbf{E}_{i,:}$  is the  $i$ -th row.

Its stability metrics are:

$$\text{Forward Stability: } \max_i \|\mathbf{E}_i\|_{\text{RMS}} = \Theta(1) \quad (4)$$

$$\text{Dependency Stability: } \max_{i,j} \|\mathbf{E}_i - \mathbf{E}_j\|_{\text{RMS}} = \Theta(1) \quad (5)$$

$$\text{Update Stability: } \max_i \|\Delta\mathbf{E}_i\|_{\text{RMS}} = \Theta(1) \quad (6)$$

Since  $\max_{i,j} \|\mathbf{E}_i - \mathbf{E}_j\|_{\text{RMS}} \leq 2 \max_i \|\mathbf{E}_i\|_{\text{RMS}}$ , these essentially measure the max row norm (scaled by  $1/\sqrt{d}$ ) of  $\mathbf{E}$  or  $\Delta\mathbf{E}$ .

To find the steepest descent direction for Embedding layers, we solve:

$$\min_{\Delta\mathbf{E}} \langle \mathbf{G}, \Delta\mathbf{E} \rangle \quad \text{s.t.} \quad \max_i \underbrace{\|\Delta\mathbf{E}_i\|_{\text{RMS}}}_{\|\Delta\mathbf{E}_i\|_2/\sqrt{d}} \leq \eta \quad (7)$$

Using Cauchy-Schwarz inequality, we derive:

$$\langle \mathbf{G}, \Delta\mathbf{E} \rangle = \sum_{i=1}^{|V|} \langle \mathbf{G}_i, \Delta\mathbf{E}_i \rangle \geq - \sum_{i=1}^{|V|} \|\mathbf{G}_i\|_2 \times \|\Delta\mathbf{E}_i\|_2 \geq -\eta\sqrt{d} \sum_{i=1}^{|V|} \|\mathbf{G}_i\|_2 \quad (8)$$

Equality holds when  $\Delta\mathbf{E}_i = -\eta\mathbf{G}_i/\|\mathbf{G}_i\|_{\text{RMS}}$ , meaning the appropriate steepest descent for Embedding layers is Row-wise RMS Normalized SGD.

### 3 LM Heads

Now consider LM Heads. Despite appearing as linear layers (input  $\mathbf{x} \in \mathbb{R}^d$ , weight  $\mathbf{W} \in \mathbb{R}^{d \times |V|}$ , output  $\mathbf{x}\mathbf{W} \in \mathbb{R}^{|V|}$ ), they are not suitable for Muon.

#### 3.1 Responsibility to Loss

The reason is that LM Heads are responsible for computing the loss during training:

$$\ell(\mathbf{x}, t; \mathbf{W}) = \log \sum_{i=1}^{|V|} e^{\langle \mathbf{x}, \boldsymbol{\omega}_i \rangle} - \langle \mathbf{x}, \boldsymbol{\omega}_t \rangle = \log \sum_{i=1}^{|V|} e^{\langle \mathbf{x}, \boldsymbol{\omega}_i - \boldsymbol{\omega}_t \rangle} \quad (9)$$

where  $\boldsymbol{\omega}_i \triangleq \mathbf{W}_{:,i}$  is the  $i$ -th column of  $\mathbf{W}$ .

The three stability metrics cannot be computed exactly due to the complex nonlinearity, but we derive tight upper bounds.

#### 3.2 Forward Stability

$$\begin{aligned} \ell(\mathbf{x}, t; \mathbf{W}) &= \log \sum_{i=1}^{|V|} e^{\langle \mathbf{x}, \boldsymbol{\omega}_i - \boldsymbol{\omega}_t \rangle} \leq \log \left( |V| \max_i e^{\langle \mathbf{x}, \boldsymbol{\omega}_i - \boldsymbol{\omega}_t \rangle} \right) \\ &= \log |V| + \max_i \langle \mathbf{x}, \boldsymbol{\omega}_i - \boldsymbol{\omega}_t \rangle \\ &\leq \log |V| + \max_i \|\mathbf{x}\|_2 \|\boldsymbol{\omega}_i - \boldsymbol{\omega}_t\|_2 \end{aligned} \quad (10)$$

$$\begin{aligned} \text{Forward Stability:} \quad \max_{t, \|\mathbf{x}\|_{\text{RMS}}=1} \ell(\mathbf{x}, t; \mathbf{W}) &\leq \log |V| + d \max_{i,t} \|\boldsymbol{\omega}_i - \boldsymbol{\omega}_t\|_{\text{RMS}} \\ &\leq \log |V| + 2d \max_i \|\boldsymbol{\omega}_i\|_{\text{RMS}} \end{aligned} \quad (11)$$

To ensure  $\Theta(1)$ , the initialization variance should be  $\Theta(1/d^2)$ .

#### 3.3 Key Inequality

We prove:

$$\left| \log \sum_{i=1}^n e^{a_i} - \log \sum_{i=1}^n e^{b_i} \right| \leq \max_i |a_i - b_i| \quad (12)$$

Using this and Cauchy-Schwarz, we derive bounds for the other metrics.

### 3.4 Dependency Stability

$$\begin{aligned}
|\ell(\mathbf{x}_1, t_1; \mathbf{W}) - \ell(\mathbf{x}_2, t_2; \mathbf{W})| &\leq \max_i |\langle \mathbf{x}_1, \boldsymbol{\omega}_i - \boldsymbol{\omega}_{t_1} \rangle - \langle \mathbf{x}_2, \boldsymbol{\omega}_i - \boldsymbol{\omega}_{t_2} \rangle| \\
&\leq \max_i (\|\mathbf{x}_1\|_2 \|\boldsymbol{\omega}_i - \boldsymbol{\omega}_{t_1}\|_2 + \|\mathbf{x}_2\|_2 \|\boldsymbol{\omega}_i - \boldsymbol{\omega}_{t_2}\|_2) \\
&= d \max_i (\|\mathbf{x}_1\|_{\text{RMS}} \|\boldsymbol{\omega}_i - \boldsymbol{\omega}_{t_1}\|_{\text{RMS}} + \|\mathbf{x}_2\|_{\text{RMS}} \|\boldsymbol{\omega}_i - \boldsymbol{\omega}_{t_2}\|_{\text{RMS}})
\end{aligned} \tag{13}$$

$$\begin{aligned}
\text{Dependency Stability: } \max_{t_1, t_2} |\ell(\mathbf{x}_1, t_1; \mathbf{W}) - \ell(\mathbf{x}_2, t_2; \mathbf{W})| &\leq d \max_{i, t_1, t_2} (\|\boldsymbol{\omega}_i - \boldsymbol{\omega}_{t_1}\|_{\text{RMS}} + \|\boldsymbol{\omega}_i - \boldsymbol{\omega}_{t_2}\|_{\text{RMS}}) \\
\|\mathbf{x}_1\|_{\text{RMS}} &= 1 \\
\|\mathbf{x}_2\|_{\text{RMS}} &= 1 \\
&\leq 4d \max_i \|\boldsymbol{\omega}_i\|_{\text{RMS}}
\end{aligned} \tag{14}$$

### 3.5 Update Stability

$$\begin{aligned}
|\ell(\mathbf{x}, t; \mathbf{W} + \Delta \mathbf{W}) - \ell(\mathbf{x}, t; \mathbf{W})| &\leq \max_i |\langle \mathbf{x}, \boldsymbol{\omega}_i + \Delta \boldsymbol{\omega}_i - \boldsymbol{\omega}_t - \Delta \boldsymbol{\omega}_t \rangle - \langle \mathbf{x}, \boldsymbol{\omega}_i - \boldsymbol{\omega}_t \rangle| \\
&= \max_i |\langle \mathbf{x}, \Delta \boldsymbol{\omega}_i - \Delta \boldsymbol{\omega}_t \rangle| \\
&\leq d \max_i \|\mathbf{x}\|_{\text{RMS}} \|\Delta \boldsymbol{\omega}_i - \Delta \boldsymbol{\omega}_t\|_{\text{RMS}}
\end{aligned} \tag{15}$$

$$\begin{aligned}
\text{Update Stability: } \max_{t, \|\mathbf{x}\|_{\text{RMS}}=1} |\ell(\mathbf{x}, t; \mathbf{W} + \Delta \mathbf{W}) - \ell(\mathbf{x}_2, t_2; \mathbf{W})| &\leq d \max_{i, t} \|\Delta \boldsymbol{\omega}_i - \Delta \boldsymbol{\omega}_t\|_{\text{RMS}} \\
&\leq 2d \max_i \|\Delta \boldsymbol{\omega}_i\|_{\text{RMS}}
\end{aligned} \tag{16}$$

The steepest descent for LM Heads is also Normalized SGD, but applied column-wise. Initialization standard deviation and learning rate scale as  $\Theta(1/d)$ , unlike the  $\Theta(1)$  scaling for Embedding layers.

## 4 Other Modules

### 4.1 Hadamard Product

After RMS Norm, we often multiply by a  $\boldsymbol{\gamma}$  vector:  $(\mathbf{x}/\|\mathbf{x}\|_{\text{RMS}}) \odot \boldsymbol{\gamma}$ . This is equivalent to matrix multiplication with a diagonal matrix  $\text{diag}(\boldsymbol{\gamma})$ , which can be analyzed as a special linear layer.

Initialization should be identity matrix (i.e.,  $\boldsymbol{\gamma}$  initialized to all ones), and the optimizer becomes SignSGD.

### 4.2 Linear Bias Terms

For linear layers with bias  $\mathbf{b}$ , the stability metrics become:

$$\text{Forward Stability: } \max_{\|\mathbf{x}\|_{\text{RMS}}=1} \|\mathbf{x}\mathbf{W} + \mathbf{b}\|_{\text{RMS}} \tag{17}$$

$$\text{Dependency Stability: } \max_{\|\mathbf{x}_1\|_{\text{RMS}}=\|\mathbf{x}_2\|_{\text{RMS}}=1} \|\mathbf{x}_1\mathbf{W} - \mathbf{x}_2\mathbf{W}\|_{\text{RMS}} \tag{18}$$

$$\text{Update Stability: } \max_{\|\mathbf{x}\|_{\text{RMS}}=1} \|\mathbf{x}\Delta \mathbf{W} + \Delta \mathbf{b}\|_{\text{RMS}} \tag{19}$$

Bias is typically initialized to zero and optimized using Normalized SGD.

### 4.3 Attention Scaling

For attention mechanisms with  $\mathbf{q} = \mathbf{x}\mathbf{W}_q, \mathbf{k} = \mathbf{x}\mathbf{W}_k$ , we have:

$$|\langle \mathbf{q}, \mathbf{k} \rangle| \leq \|\mathbf{q}\|_2 \|\mathbf{k}\|_2 = d \|\mathbf{q}\|_{\text{RMS}} \|\mathbf{k}\|_{\text{RMS}} \quad (20)$$

This suggests a scaling factor of  $\Theta(1/d)$ , which complements the common  $1/\sqrt{d}$  initialization scaling.

## 5 Summary

The key results are summarized below:

	Input	Parameters	Output	Initial Variance	Steepest Descent
Linear	$\mathbf{x}$	$\mathbf{W} \in \mathbb{R}^{d_{in} \times d_{out}}$ $\mathbf{b} \in \mathbb{R}^{d_{out}}$	$\mathbf{x}\mathbf{W} + \mathbf{b}$	$\mathbf{W}: \sqrt{\frac{d_{out}}{d_{in}}} \frac{1}{\sqrt{d_{in} + \sqrt{d_{out}}}}$ $\mathbf{b}: 0$	$\Delta \mathbf{W} = -\eta \sqrt{\frac{d_{out}}{d_{in}}} \text{msign}(\mathbf{G})$ $\Delta \mathbf{b} = -\eta \frac{\mathbf{g}}{\ \mathbf{g}\ _{\text{RMS}}}$
Embedding	$i$	$\mathbf{E} \in \mathbb{R}^{ V  \times d}$	$\mathbf{E}_{i,:}$	1	$\Delta \mathbf{E}_{i,:} = -\eta \frac{\mathbf{G}_{i,:}}{\ \mathbf{G}_{i,:}\ _{\text{RMS}}}$
LM Head	$\mathbf{x}, t$	$\mathbf{W} \in \mathbb{R}^{d \times  V }$	$\log \sum_{i=1}^{ V } e^{\langle \mathbf{x}, \mathbf{W}_{:,i} - \mathbf{W}_{:,t} \rangle}$	$\frac{1}{d^2}$	$\Delta \mathbf{W}_{:,i} = -\frac{\eta}{d} \frac{\mathbf{G}_{:,i}}{\ \mathbf{G}_{:,i}\ _{\text{RMS}}}$
RMS Norm	$\mathbf{x}$	$\gamma \in \mathbb{R}^d$	$\frac{\mathbf{x}}{\ \mathbf{x}\ _{\text{RMS}}} \odot \gamma$	1	$\Delta \gamma = -\eta \text{sign}(\mathbf{g})$

This shows that stability metrics provide a unified framework for analyzing and optimizing different types of layers, and explain why certain components require special treatment.