

Controlled LLM Training On The Spectral Sphere

Spectral Sphere Optimizer

Author

*Tian Xie, Haoming Luo, Haoyu Tang,
Yiwen Hu, J.K. Liu, Qingnan Ren,
Yang Wang, W.X. Zhao, Rui Yan,
Bing Su, Chong Luo, Baining Guo*

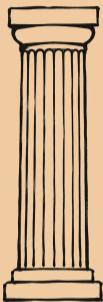


GitHub



WeChat

Two Pillars



Convergence Speed

How fast does loss drop?
A: Steepest descent under
a good norm



Measurable Stability

Do activations explode?
A: μP condition for
feature learning

LLM training is, at its core, a pursuit of **convergence speed** grounded in the necessity of **stability**.

Fantastic Optimizers



Each one is steepest descent
under a different geometry.

**SSO: a fast & stable
optimizer under spectral sphere**

Steepest Descent

1. Linearize the loss
2. cap the step under a norm
3. minimize.

$$\Delta \mathbf{W} := \arg \min_{\|\Delta \mathbf{W}\| \leq \eta} \underbrace{\{ \mathcal{L}(\mathbf{W}) + \langle \mathbf{G}, \Delta \mathbf{W} \rangle \}}_{\text{first-order model}}$$

$$\Rightarrow \Delta \mathbf{W} = -\eta \cdot \arg \max_{\|\mathbf{T}\| \leq 1} \langle \mathbf{G}, \mathbf{T} \rangle$$

Frobenius $\|\cdot\|_F$

$$\Delta \mathbf{W} = -\eta \mathbf{G}$$

SGD

vector ℓ_∞

$$\Delta \mathbf{W} = -\eta \text{sign}(\mathbf{G})$$

Adam (no EMA)

matrix spectral $\|\cdot\|_2$

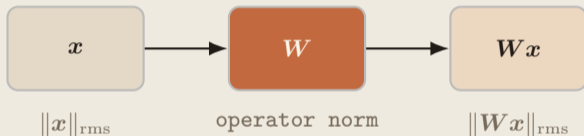
$$\mathbf{G} = \mathbf{U}\Sigma\mathbf{V}^\top \mapsto \mathbf{U}\mathbf{V}^\top$$

Muon (msign)

Muon- steepest descent under the *spectral norm*

μ P Condition: Hidden Stability

Stable training = activations stay $\Theta(1)$. RMS norm $\|x\|_{\text{rms}} = \|x\|_2/\sqrt{d}$.



Forward stability — a layer must preserve the activation scale :

$$\|W\|_{\text{rms} \rightarrow \text{rms}} := \max_{x \neq 0} \frac{\|Wx\|_{\text{rms}}}{\|x\|_{\text{rms}}} = \Theta(1)$$

μ P Condition: Hidden Stability

Plug in $\|\cdot\|_{\text{rms}} = \|\cdot\|_2/\sqrt{d}$: the \sqrt{d} 's pull out, leaving the ordinary spectral norm $\|\mathbf{W}\|_2$ scaled by the widths $d_{\text{in}} \rightarrow d_{\text{out}}$.

$$\|\mathbf{W}\|_{\text{rms} \rightarrow \text{rms}} = \sqrt{\frac{d_{\text{in}}}{d_{\text{out}}}} \max_{\mathbf{x} \neq \mathbf{0}} \frac{\|\mathbf{W}\mathbf{x}\|_2}{\|\mathbf{x}\|_2} = \sqrt{\frac{d_{\text{in}}}{d_{\text{out}}}} \|\mathbf{W}\|_2 = \Theta(1)$$

Forward \rightarrow **weight**

$$\|\mathbf{W}\|_2 = \Theta(\sqrt{d_{\text{out}}/d_{\text{in}}})$$

Update \rightarrow **optimizer**

$$\|\Delta\mathbf{W}\|_2 = \Theta(\sqrt{d_{\text{out}}/d_{\text{in}}})$$

Spectral μ P Condition: $\|\mathbf{W}\|_2 = \|\Delta\mathbf{W}\|_2 = \Theta(\sqrt{d_{\text{out}}/d_{\text{in}}}) =: R$

Muon: Half-Aligned μP

Muon controls $\|\Delta W\|_2$ ✓

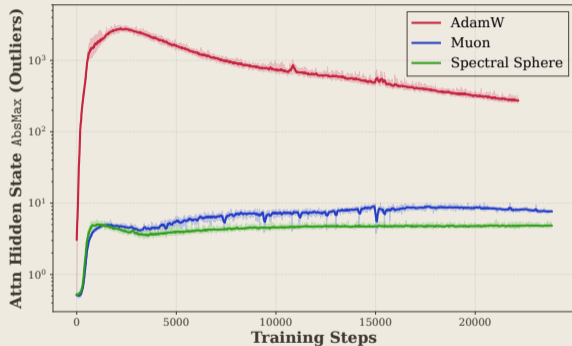
Muon does *not* control $\|W\|_2$ ✗

→ effective step $\|\Delta W\|_2/\|W\|_2$ drifts

→ attention-logit blow-ups in long runs

→ needs *non-essential*

SandwichNorm / QK-Norm / softcap patches



Attention activation AbsMax, Dense-1.7B pre-training

Spectral Sphere Optimizer

What if one optimizer could simultaneously satisfy the steepest-descent property for **convergence speed**, and the strict μP constraints for **fundamental stability**?

The Objective

Pin **both** the weight \mathbf{W} *and* its update onto the spectral sphere of radius $R = \Theta(\sqrt{d_{\text{out}}/d_{\text{in}}})$.
Parameterize $\Delta\mathbf{W} = \eta R \Phi$ and solve for the unit step Φ :

$$\max_{\Phi} \langle \mathbf{G}, \Phi \rangle$$

The Objective

Pin **both** the weight W and its update onto the spectral sphere of radius $R = \Theta(\sqrt{d_{\text{out}}/d_{\text{in}}})$.
Parameterize $\Delta W = \eta R \Phi$ and solve for the unit step Φ :

$$\max_{\Phi} \langle G, \Phi \rangle$$

subject to

$$\|\Phi\|_2 = 1$$

steepest descent under the spectral norm

The Objective

Pin **both** the weight W and its update onto the spectral sphere of radius $R = \Theta(\sqrt{d_{\text{out}}/d_{\text{in}}})$.
Parameterize $\Delta W = \eta R \Phi$ and solve for the unit step Φ :

$$\max_{\Phi} \langle G, \Phi \rangle$$

subject to

$$\|\Phi\|_2 = 1$$

steepest descent under the spectral norm

$$\|W - \eta R \Phi\|_2 = \|W\|_2 = R$$

W stays on the sphere of radius R

First Order: From Sphere to Tangent

Spectral norm is differentiable when the top singular value is simple (a.s. for random matrices). Its gradient is the **top singular pair**:

$$\Theta := \nabla_{\mathbf{W}} \|\mathbf{W}\|_2 = \mathbf{u}_1 \mathbf{v}_1^\top.$$

First Order: From Sphere to Tangent

Spectral norm is differentiable when the top singular value is simple (a.s. for random matrices). Its gradient is the **top singular pair**:

$$\Theta := \nabla_{\mathbf{W}} \|\mathbf{W}\|_2 = \mathbf{u}_1 \mathbf{v}_1^\top.$$

First-order Taylor expansion of the moving-weight constraint about \mathbf{W} :

$$\|\mathbf{W} - \eta R \Phi\|_2 = \|\mathbf{W}\|_2 - \eta R \langle \Theta, \Phi \rangle + \mathcal{O}(\eta^2 R^2).$$

First Order: From Sphere to Tangent

Spectral norm is differentiable when the top singular value is simple (a.s. for random matrices). Its gradient is the **top singular pair**:

$$\Theta := \nabla_{\mathbf{W}} \|\mathbf{W}\|_2 = \mathbf{u}_1 \mathbf{v}_1^\top.$$

First-order Taylor expansion of the moving-weight constraint about \mathbf{W} :

$$\|\mathbf{W} - \eta R \Phi\|_2 = \|\mathbf{W}\|_2 - \eta R \langle \Theta, \Phi \rangle + \mathcal{O}(\eta^2 R^2).$$

Matching $\|\mathbf{W} - \eta R \Phi\|_2 = \|\mathbf{W}\|_2$ at first order \Rightarrow **tangent constraint** $\langle \Theta, \Phi \rangle = 0$.

First Order: From Sphere to Tangent

Spectral norm is differentiable when the top singular value is simple (a.s. for random matrices). Its gradient is the **top singular pair**:

$$\Theta := \nabla_{\mathbf{W}} \|\mathbf{W}\|_2 = \mathbf{u}_1 \mathbf{v}_1^\top.$$

First-order Taylor expansion of the moving-weight constraint about \mathbf{W} :

$$\|\mathbf{W} - \eta R \Phi\|_2 = \|\mathbf{W}\|_2 - \eta R \langle \Theta, \Phi \rangle + \mathcal{O}(\eta^2 R^2).$$

Matching $\|\mathbf{W} - \eta R \Phi\|_2 = \|\mathbf{W}\|_2$ at first order \Rightarrow **tangent constraint** $\langle \Theta, \Phi \rangle = 0$.

the problem collapses to

$$\max_{\Phi} \langle \mathbf{G}, \Phi \rangle \quad \text{s.t.} \quad \|\Phi\|_2 = 1, \quad \langle \Theta, \Phi \rangle = 0$$

The Solution: One Rank-1 Correction

A single **Lagrange multiplier** λ folds the tangent constraint into the objective:

$$\mathcal{L}(\Phi; \lambda) = \langle \mathbf{G} + \lambda \Theta, \Phi \rangle, \quad \|\Phi\|_2 = 1.$$

The Solution: One Rank-1 Correction

A single **Lagrange multiplier** λ folds the tangent constraint into the objective:

$$\mathcal{L}(\Phi; \lambda) = \langle \mathbf{G} + \lambda \Theta, \Phi \rangle, \quad \|\Phi\|_2 = 1.$$

Steepest descent under the spectral norm \Rightarrow optimum is an `msign`:

$$\Phi^*(\lambda) = \text{msign}(\mathbf{G} + \lambda \Theta)$$

The Solution: One Rank-1 Correction

A single **Lagrange multiplier** λ folds the tangent constraint into the objective:

$$\mathcal{L}(\Phi; \lambda) = \langle \mathbf{G} + \lambda \Theta, \Phi \rangle, \quad \|\Phi\|_2 = 1.$$

Steepest descent under the spectral norm \Rightarrow optimum is an `msign`:

$$\Phi^*(\lambda) = \text{msign}(\mathbf{G} + \lambda \Theta)$$

pick λ^* so the tangent constraint holds

$$h(\lambda) := \langle \Theta, \text{msign}(\mathbf{G} + \lambda \Theta) \rangle = 0$$

The Solution: One Rank-1 Correction

A single **Lagrange multiplier** λ folds the tangent constraint into the objective:

$$\mathcal{L}(\Phi; \lambda) = \langle \mathbf{G} + \lambda \Theta, \Phi \rangle, \quad \|\Phi\|_2 = 1.$$

Steepest descent under the spectral norm \Rightarrow optimum is an `msign`:

$$\Phi^*(\lambda) = \text{msign}(\mathbf{G} + \lambda \Theta)$$

pick λ^* so the tangent constraint holds

$$h(\lambda) := \langle \Theta, \text{msign}(\mathbf{G} + \lambda \Theta) \rangle = 0$$

$\lambda=0 \Rightarrow$ exactly Muon's `msign`(\mathbf{G}).

SSO adds *one rank-1 correction* $\lambda \Theta$ that pins the step to the sphere.

Solving for λ^*

Monotone. $h(\lambda)$ is non-decreasing, sweeping $-1 \rightarrow +1$ as λ runs over \mathbb{R} .

Solving for λ^*

Monotone. $h(\lambda)$ is non-decreasing, sweeping $-1 \rightarrow +1$ as λ runs over \mathbb{R} .

Localized. $|\lambda^*| \leq 2\|G\|_*$ — a finite search window.

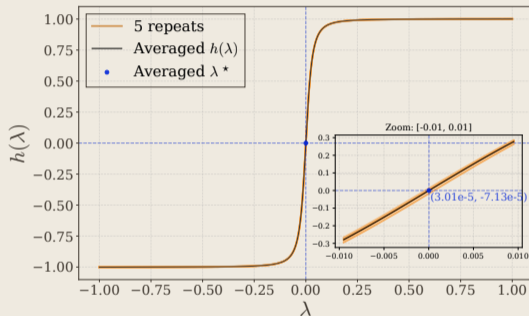
Solving for λ^*

Monotone. $h(\lambda)$ is non-decreasing, sweeping $-1 \rightarrow +1$ as λ runs over \mathbb{R} .

Localized. $|\lambda^*| \leq 2\|G\|_*$ — a finite search window.

Bracket from $\lambda=0$ (expand against sign $h(0)$), then **bisect** to tolerance.

a 1-D root find -- a handful of `msign` calls per step.



$h(\lambda)$ on a random 1024×3072 matrix; root λ^* sits close to 0

Second Order: Retraction

The $\mathcal{O}(\eta^2 R^2)$ remainder accumulates and drifts \mathbf{W} off the sphere. Restore $\|\mathbf{W}\|_2 = R$ *exactly* with a one-line **retraction**:

$$\mathbf{W} \leftarrow \mathbf{W} \cdot \frac{R}{\|\mathbf{W}\|_2}$$

Second Order: Retraction

The $\mathcal{O}(\eta^2 R^2)$ remainder accumulates and drifts \mathbf{W} off the sphere. Restore $\|\mathbf{W}\|_2 = R$ *exactly* with a one-line **retraction**:

$$\mathbf{W} \leftarrow \mathbf{W} \cdot \frac{R}{\|\mathbf{W}\|_2}$$

Free. A pre-update step reuses *one* power iteration for $\|\mathbf{W}\|_2$ (cached power iteration).

Bounded. $\|\mathbf{W}\|_2=R \Rightarrow \|\mathbf{W}\|_F \leq \sqrt{\text{rank}} R.$

Second Order: Retraction

The $\mathcal{O}(\eta^2 R^2)$ remainder accumulates and drifts \mathbf{W} off the sphere. Restore $\|\mathbf{W}\|_2 = R$ *exactly* with a one-line **retraction**:

$$\mathbf{W} \leftarrow \mathbf{W} \cdot \frac{R}{\|\mathbf{W}\|_2}$$

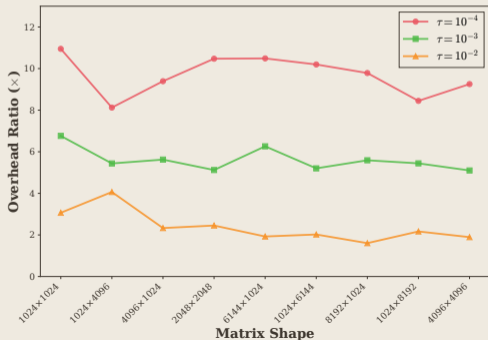
Free. A pre-update step reuses *one* power iteration for $\|\mathbf{W}\|_2$ (cached power iteration).

Bounded. $\|\mathbf{W}\|_2=R \Rightarrow \|\mathbf{W}\|_F \leq \sqrt{\text{rank}} R.$

No weight decay on hidden 2D weights. One fewer hyperparameter to tune, good news for scaling law fitting.

Closing the Infra Cost Gap

Stage	Time (ms)	Δ vs prev.
Naive baseline (no opt.)	10 928	—
+ load balance & All-Gather	9 366	-14.3%
+ adaptive kernel & multi-stream	9 383	-8.8%
+ BF16 & torch.compile (final)	7 666	-18.3%



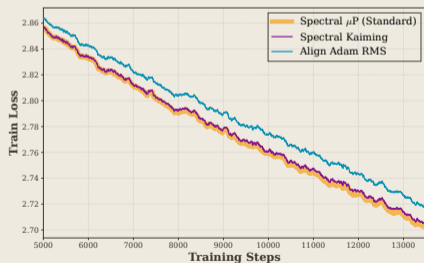
AdamW 6.73s · Muon 6.88s
MuonSphere 6.95s · SSO 7.67s

For tight infra budgets, MuonSphere ($\lambda=0$) keeps the activation control at near-zero overhead.

The Practical Recipe: LR Scaler

Unified view: each LR scaler keeps a consistent **effective step size** under a chosen **norm metric** and **initialization scheme**.

$$\frac{\|\Delta \mathbf{W}\|_{\mathcal{M}}}{\|\mathbf{W}\|_{\mathcal{M}}} = \frac{\|\eta R \Phi\|_{\mathcal{M}}}{\|\mathbf{W}\|_{\mathcal{M}}} = \eta$$



Spectral-manifold optimization needs a scaler explicitly calibrated to the spectral norm

Spectral μP

$$R = \sqrt{d_{\text{out}}/d_{\text{in}}}$$

spectral norm; μP init

Align-Adam-RMS

$$R = 0.2\sqrt{\max(d_{\text{out}}, d_{\text{in}})}$$

l2 norm; standard init

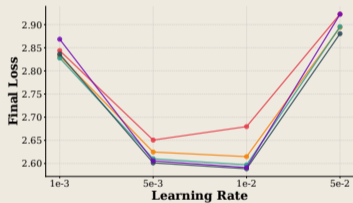
Spectral Kaiming

$$R = \sqrt{\max(1, d_{\text{out}}/d_{\text{in}})}$$

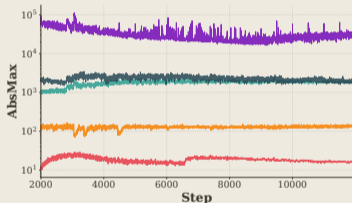
spectral norm; Kaiming init

The Practical Recipe: Spectral Radius

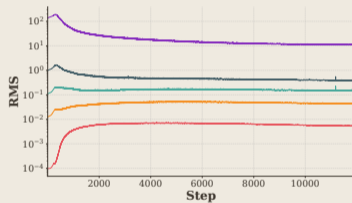
$$R = c\sqrt{d_{\text{out}}/d_{\text{in}}} \quad c \text{ sets branch output magnitude relative to the residual stream.}$$



optimization: moderate c wins,
 ~ 2.0 is best



FFN AbsMax: follows the radius
scale

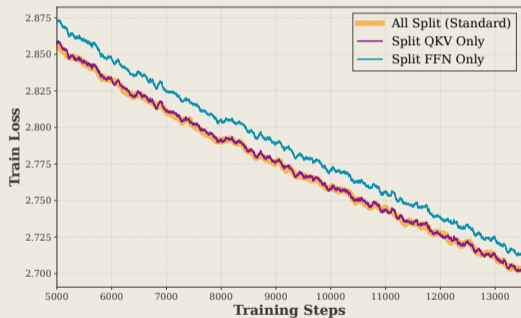


FFN RMS: clear power-law response
to c

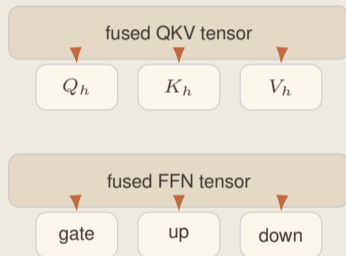
Tune c to control the signal-to-noise ratio along the deep residual path: attention/FFN branches become **directly scale-controlled** instead of drifting freely.

The Practical Recipe: Module Granularity

Megatron may fuse QKV or SwiGLU matrices into one tensor, but those slices do not play one functional role. SSO treats each slice as an **atomic function unit**.



splitting QKV per head gives the dominant gain



Best granularity: split attention QKV per head; separate FFN gate/up.

Empirical Evidence

100B-token pre-training on OLMo-Mix-1124

deterministic data order, global batch $\sim 4\text{M}$ tokens, BF16 training

Scaling

70M \rightarrow 1.8B

LR transfer

Dense

1.8B

speed

MoE

8B-A1B

routing

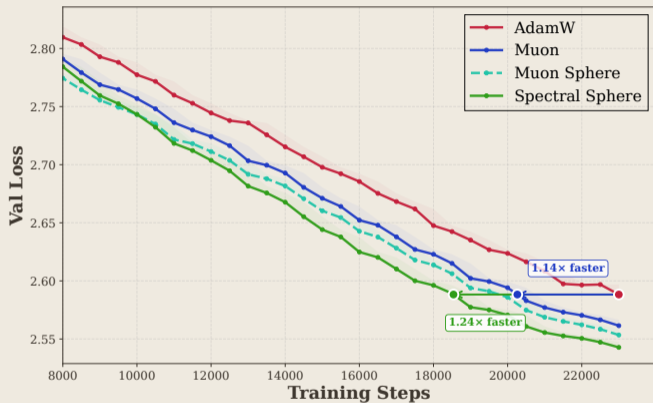
DeepNet

200 layers

depth

We test the optimizer where instability usually appears: width scaling, activation magnitude, MoE routing, and very deep residual paths.

Dense 1.8B: Faster Loss Descent



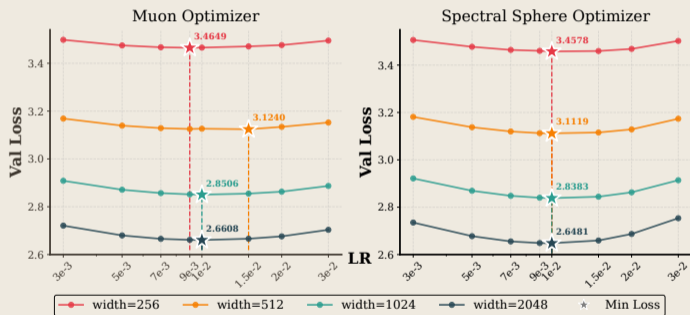
validation loss on 100B-token Dense 1.7B training

Reference point. AdamW reaches validation loss 2.588 at 23k steps.

Same loss sooner. Muon: 12% fewer steps; **SSO:** 19% fewer steps.

SSO gives a 1.24× speed-up over AdamW, even in a setup tuned for AdamW.

μP Width Transfer



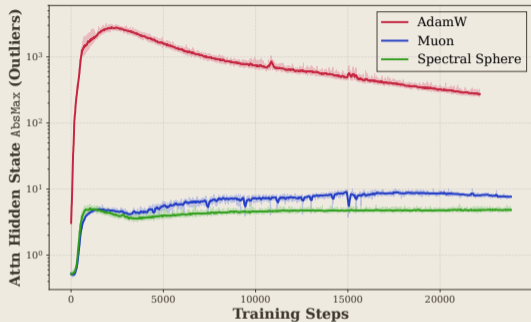
learning-rate sweep across $25\times$ model sizes: 70M \rightarrow 1.8B

Target behavior. μP should make the optimal LR transfer across width.

Muon is not enough. It controls updates, but its optimal LR still drifts with width.

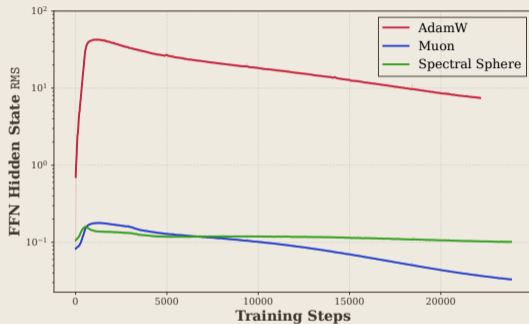
SSO transfers. Stable optimal LR, with lower optimal loss than Muon.

Hidden Activations: Bounded Dynamics



attention AbsMax: outlier magnitude

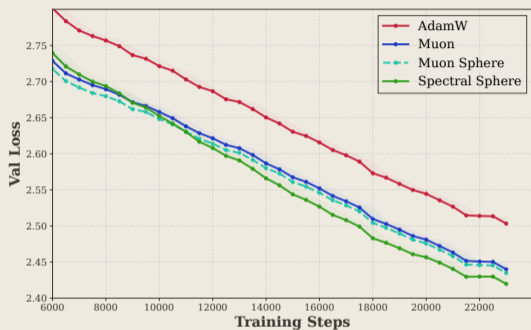
SSO keeps activation RMS at $\Theta(1)$ scale throughout training.



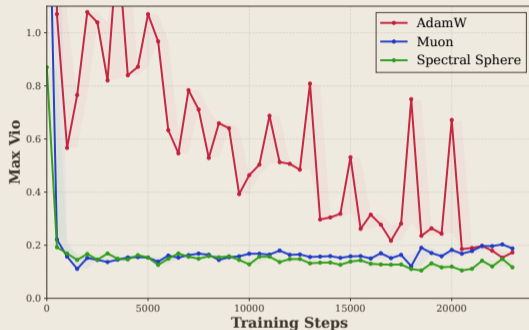
FFN RMS: activation scale

AdamW reaches $\sim 100\times$ larger activation magnitudes; **Muon** still drifts mildly.

MoE 8B-A1B: Better Router Balance



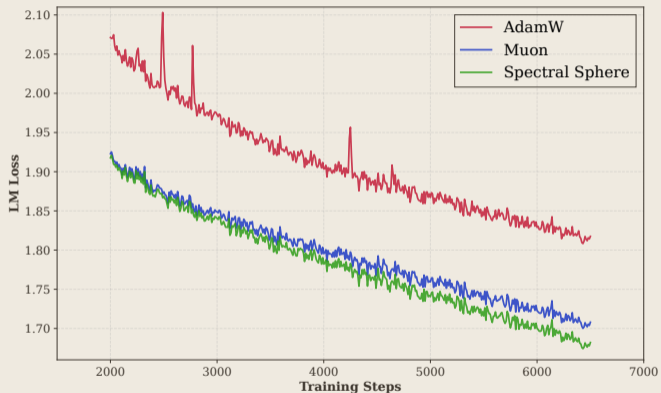
MoE validation loss



router MaxVio; lower means better load balance

Spectral normalization promotes balanced expert utilization: **SSO** has the lowest validation loss and the smallest routing violation.

DeepNet 200L: Depth Stress Test



training loss after extending the 28-layer baseline to 200 layers

Why this test. Very deep residual paths expose slow activation drift and loss spikes.

Spectral methods. Lower loss and smoother dynamics than AdamW.

SSO. Best stability and lowest loss under the 200-layer stress test.

Why Spectral Sphere?

Stiefel

all $\sigma_i = 1$
too rigid

no internal spectrum

Hyperball

$\|\mathbf{W}\|_F \leq C$
too loose

one σ can dominate

Spectral sphere

$\sigma_{\max} = R$
a balance?

worst-case bound, free interior

SSO bounds only σ_{\max} — the worst-case activation bound $\mu\mathbf{P}$ needs — and lets the interior spectrum evolve freely.

Take-Away

**SSO Steepest Descent on the Spectral Sphere
= Muon + Tangent Correction + Weight Retraction.**

✓ **Speed:** 1.24× AdamW on Dense 1.7B

✓ **Stability:** bounded activations

✓ **Practicality:** drop-in for Megatron

Acknowledgements

Many ideas in this talk were inspired by Jianlin Su's writing at kexue.fm.

Recommended reading from kexue.fm

- 1 Steepest Descent on Manifolds: 1. SGD + Hypersphere
kexue.fm/archives/11196
- 2 Steepest Descent on Manifolds: 2. Muon + Orthogonality
kexue.fm/archives/11215
- 3 Steepest Descent on Manifolds: 3. Muon + Stiefel
kexue.fm/archives/11221
- 4 Steepest Descent on Manifolds: 4. Muon + Spectral Sphere
kexue.fm/archives/11241
- 5 Beyond MuP: 4. Holding the Line on Parameter Stability
kexue.fm/archives/11729



You think you're smarter than Su-shen?

kexue.fm is all you need.

Q&A

Any questions?



GitHub



WeChat