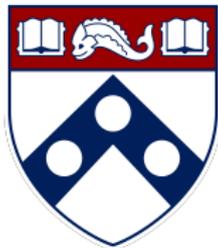


# Preconditioning benefits of spectral orthogonalization in Muon



Yuxin Chen

Wharton Statistics & Data Science

# Coauthors

---



Jianhao Ma  
UPenn



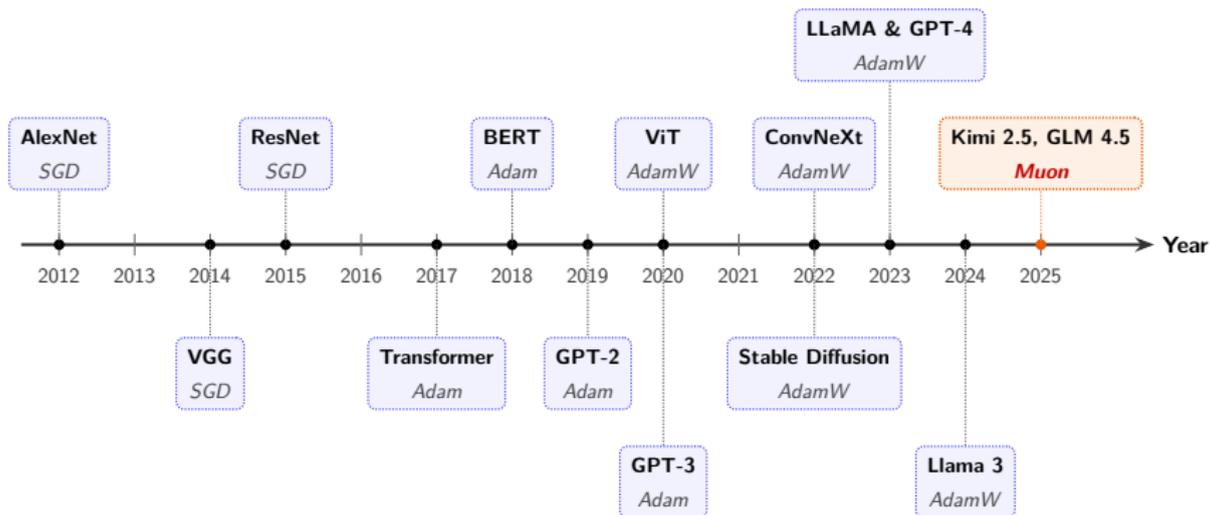
Yu Huang  
UPenn



Yuejie Chi  
Yale

“Preconditioning benefits of spectral orthogonalization in Muon,” J. Ma, Y. Huang,  
Y. Chi, Y. Chen, arXiv:2601.13474, 2026

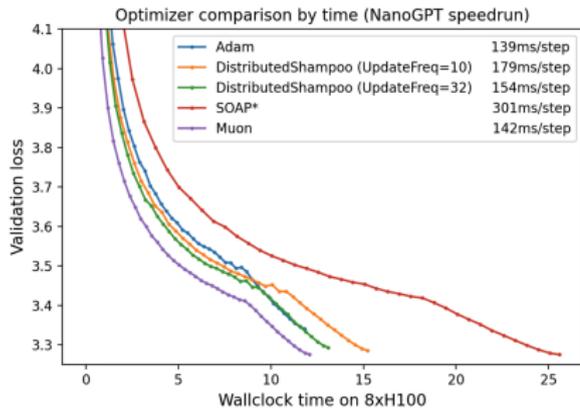
# Evolution of deep learning optimizers



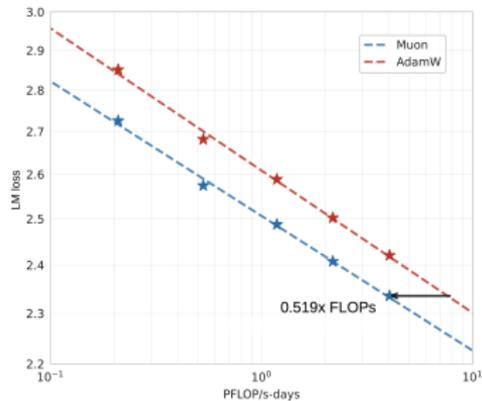
*Muon (MomentUm Orthogonalized by Newton-Schulz) optimizer*  
— Keller Jordan blog '24

## Companies are experimenting with Muon for LLM training

- Moonshot AI, Z.ai, OpenAI, DeepSeek, ...

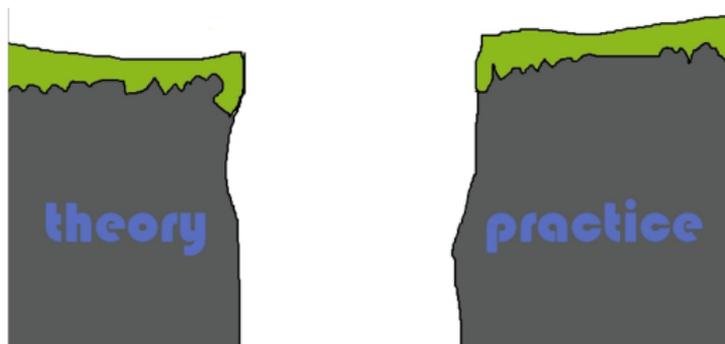


Muon's original blog

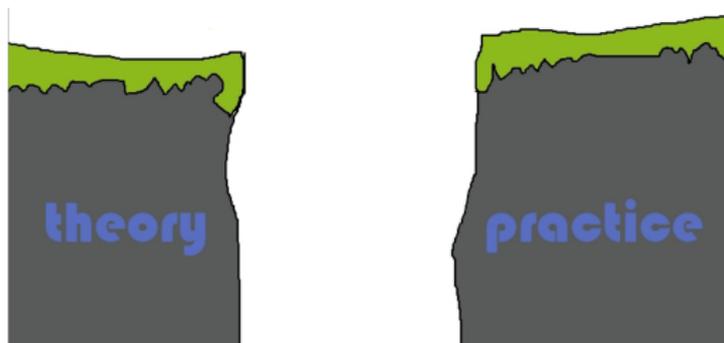


Kimi report

*Theoretically, it is mysterious why Muon performs well in practice ...*



*Theoretically, it is mysterious why Muon performs well in practice ...*



- *generic analysis*: Shen et al. '25, Li, Hong '25, Chen et al. '25, An et al. '25, Kovalev '25, ...
- *one-step analysis*: Davis, Drusvyatskiy '25, Su '25, ...
- *generalization*: Vasudeva et al. '25, Wang et al. '25, Zhang et al. '25, ...
- *earlier work*: Carlson et al. '15, Gupta et al. '18, Tuddenham et al. '22, Vyas et al. '25 ...

*Can we develop end-to-end theory to  
show provable benefits of Muon over other optimizers?*

# Muon algorithm

---

$$\text{minimize}_{\mathbf{X}} f(\mathbf{X})$$

# Muon algorithm

---

$$\text{minimize}_{\mathbf{X}} f(\mathbf{X})$$

**Muon:** for  $t = 0, 1, \dots$

$$\mathbf{B}_t = \nabla f(\mathbf{X}_t) + \mu \mathbf{B}_{t-1} \quad (\text{gradient} + \text{momentum})$$

$$\mathbf{X}_{t+1} = \mathbf{X}_t - \eta_t \text{msign}(\mathbf{B}_t) \quad (\text{spectral orthogonalization})$$

# Muon algorithm

---

$$\text{minimize}_{\mathbf{X}} \quad f(\mathbf{X})$$

**Muon:** for  $t = 0, 1, \dots$

$$\mathbf{B}_t = \nabla f(\mathbf{X}_t) + \mu \mathbf{B}_{t-1} \quad (\text{gradient} + \text{momentum})$$

$$\mathbf{X}_{t+1} = \mathbf{X}_t - \eta_t \text{msign}(\mathbf{B}_t) \quad (\text{spectral orthogonalization})$$

- $\underbrace{\text{msign}(\mathbf{Z}) := \mathbf{U}\mathbf{V}^\top}_{\text{matrix sign function}}$  if  $\mathbf{Z}$  has SVD  $\mathbf{Z} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^\top$

# Muon algorithm

---

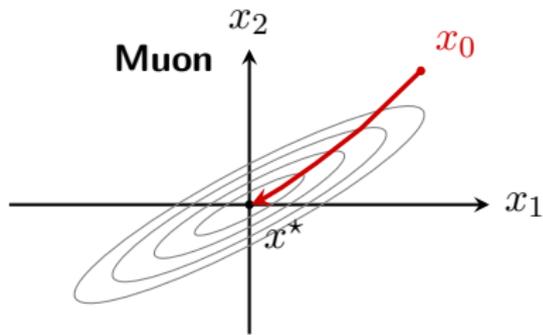
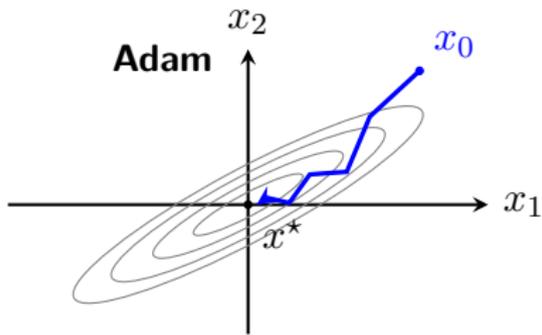
$$\text{minimize}_{\mathbf{X}} \quad f(\mathbf{X})$$

**simplified Muon (or spectral gradient method):** for  $t = 0, 1, \dots$

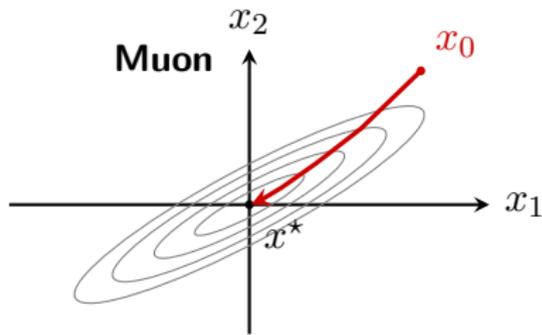
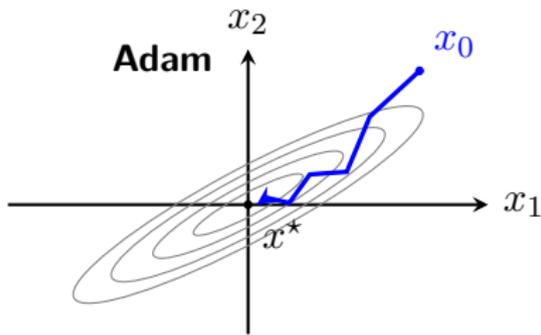
$$\mathbf{X}_{t+1} = \mathbf{X}_t - \eta_t \text{msign}(\nabla f(\mathbf{X}_t)) \quad (\text{no momentum})$$

- $\underbrace{\text{msign}(\mathbf{Z}) := \mathbf{U}\mathbf{V}^\top}_{\text{matrix sign function}}$  if  $\mathbf{Z}$  has SVD  $\mathbf{Z} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^\top$

***This talk:*** *end-to-end theory*  *of simplified Muon*



***This talk:*** *end-to-end theory of simplified Muon*  
*preconditioning effect*



***This talk:*** *end-to-end theory of simplified Muon (case studies)*  
*preconditioning effect*

- **matrix factorization**
- in-context learning of linear transformers

# Matrix factorization

---

— Chi, Lu, Chen '19

$$\underset{U \in \mathbb{R}^{d \times k}}{\text{minimize}} \quad f(U) = \frac{1}{4} \|UU^\top - M^*\|_F^2$$

- $M^* \succeq \mathbf{0}$ : rank- $r$
- condition number:  $\kappa := \sigma_{\max}(M^*)/\sigma_{\min}(M^*)$

# Matrix factorization

---

— Chi, Lu, Chen '19

$$\underset{U \in \mathbb{R}^{d \times k}}{\text{minimize}} \quad f(U) = \frac{1}{4} \|UU^\top - M^*\|_F^2$$

- $M^* \succeq \mathbf{0}$ : rank- $r$
- condition number:  $\kappa := \sigma_{\max}(M^*)/\sigma_{\min}(M^*)$
- $k = r$ : exact parameterization;  $k > r$ : over-parameterization

# Matrix factorization

---

— Chi, Lu, Chen '19

$$\underset{U \in \mathbb{R}^{d \times k}}{\text{minimize}} \quad f(U) = \frac{1}{4} \|UU^\top - M^*\|_F^2$$

- $M^* \succeq 0$ : rank- $r$
- condition number:  $\kappa := \sigma_{\max}(M^*)/\sigma_{\min}(M^*)$
- $k = r$ : exact parameterization;  $k > r$ : over-parameterization

simplified Muon :  $U_{t+1} = U_t - \eta_t \text{msign}((U_t U_t^\top - M^*)U_t)$

# Matrix factorization

— Chi, Lu, Chen '19

$$\underset{U \in \mathbb{R}^{d \times k}}{\text{minimize}} \quad f(U) = \frac{1}{4} \|UU^\top - M^*\|_F^2$$

- $M^* \succeq 0$ : rank- $r$
- condition number:  $\kappa := \sigma_{\max}(M^*)/\sigma_{\min}(M^*)$
- $k = r$ : exact parameterization;  $k > r$ : over-parameterization

$$\text{simplified Muon : } U_{t+1} = U_t - \eta_t \text{msign}((U_t U_t^\top - M^*)U_t)$$

vs.

$$\text{GD : } U_{t+1} = U_t - \eta_t ((U_t U_t^\top - M^*)U_t)$$
$$\underbrace{\text{SignGD}}_{\text{simplified Adam}} : U_{t+1} = U_t - \eta_t \text{sign}((U_t U_t^\top - M^*)U_t)$$

# Convergence theory of Muon

---

simplified Muon :  $U_{t+1} = U_t - \eta_t \text{msign}((U_t U_t^\top - M^*)U_t)$

# Convergence theory of Muon

---

simplified Muon :  $U_{t+1} = U_t - \eta_t \text{msign}((U_t U_t^\top - M^*)U_t)$

## Theorem 1 (exactly parameterized ( $k = r$ ))

- *stepsize*:  $\eta_t = C_{\eta,t} \sqrt{\sigma_{\max}^*} \rho^t$  for  $2/3 \leq \rho < 1$ ,  $C_{\eta,t} \stackrel{i.i.d.}{\sim} \text{Unif}(1, 2)$
- *initialization*:  $U_0 = \alpha O$  w/ random orthonormal  $O$  and small  $\alpha$

# Convergence theory of Muon

---

simplified Muon :  $U_{t+1} = U_t - \eta_t \text{msign}((U_t U_t^\top - M^*)U_t)$

## Theorem 1 (exactly parameterized ( $k = r$ ))

- *stepsize*:  $\eta_t = C_{\eta,t} \sqrt{\sigma_{\max}^*} \rho^t$  for  $2/3 \leq \rho < 1$ ,  $C_{\eta,t} \stackrel{i.i.d.}{\sim} \text{Unif}(1, 2)$
- *initialization*:  $U_0 = \alpha O$  w/ random orthonormal  $O$  and small  $\alpha$

With prob. at least 0.99, one has  $\|U_T U_T^\top - M^*\| \leq \varepsilon$  if

$$T = \frac{1}{1 - \rho} \log \left( \frac{16\sigma_{\max}^*}{(1 - \rho)^2 \varepsilon} \right)$$

# Convergence theory of Muon

---

simplified Muon :  $U_{t+1} = U_t - \eta_t \text{msign}((U_t U_t^\top - M^*)U_t)$

## Theorem 1 (exactly parameterized ( $k = r$ ))

- *stepsize*:  $\eta_t = C_{\eta,t} \sqrt{\sigma_{\max}^*} \rho^t$  for  $2/3 \leq \rho < 1$ ,  $C_{\eta,t} \stackrel{i.i.d.}{\sim} \text{Unif}(1, 2)$
- *initialization*:  $U_0 = \alpha O$  w/ random orthonormal  $O$  and small  $\alpha$

With prob. at least 0.99, one has  $\|U_T U_T^\top - M^*\| \leq \varepsilon$  if

$$T = \frac{1}{1 - \rho} \log \left( \frac{16\sigma_{\max}^*}{(1 - \rho)^2 \varepsilon} \right)$$

- linear convergence

# Convergence theory of Muon

simplified Muon :  $U_{t+1} = U_t - \eta_t \text{msign}((U_t U_t^\top - M^*)U_t)$

## Theorem 1 (exactly parameterized ( $k = r$ ))

- *stepsize*:  $\eta_t = C_{\eta,t} \sqrt{\sigma_{\max}^*} \rho^t$  for  $2/3 \leq \rho < 1$ ,  $C_{\eta,t} \stackrel{i.i.d.}{\sim} \text{Unif}(1, 2)$
- *initialization*:  $U_0 = \alpha O$  w/ random orthonormal  $O$  and small  $\alpha$

With prob. at least 0.99, one has  $\|U_T U_T^\top - M^*\| \leq \varepsilon$  if

$$T = \frac{1}{1 - \rho} \log \left( \frac{16\sigma_{\max}^*}{(1 - \rho)^2 \varepsilon} \right)$$

- linear convergence
- iteration complexity independent of condition number  $\kappa$

# Convergence theory of Muon

simplified Muon :  $U_{t+1} = U_t - \eta_t \text{msign}((U_t U_t^\top - M^*)U_t)$

## Theorem 1 (exactly parameterized ( $k = r$ ))

- *stepsize*:  $\eta_t = C_{\eta,t} \sqrt{\sigma_{\max}^*} \rho^t$  for  $2/3 \leq \rho < 1$ ,  $C_{\eta,t} \stackrel{i.i.d.}{\sim} \text{Unif}(1, 2)$
- *initialization*:  $U_0 = \alpha O$  w/ random orthonormal  $O$  and small  $\alpha$

With prob. at least 0.99, one has  $\|U_T U_T^\top - M^*\| \leq \varepsilon$  if

$$T = \frac{1}{1 - \rho} \log \left( \frac{16\sigma_{\max}^*}{(1 - \rho)^2 \varepsilon} \right)$$

- linear convergence
- iteration complexity independent of condition number  $\kappa$
- geometrically decaying stepsizes

# Convergence theory of Muon

simplified Muon :  $U_{t+1} = U_t - \eta_t \text{msign}((U_t U_t^\top - M^*)U_t)$

## Theorem 1 (mildly over-parameterized ( $r \leq k < d$ ))

- *stepsize*:  $\eta_t = C_{\eta,t} \sqrt{\sigma_{\max}^*} \rho^t$  for  $2/3 \leq \rho < 1$ ,  $C_{\eta,t} \stackrel{i.i.d.}{\sim} \text{Unif}(1, 2)$
- *initialization*:  $U_0 = \alpha O$  w/ random orthonormal  $O$  and small  $\alpha$

With prob. at least 0.99, one has  $\|U_T U_T^\top - M^*\| \leq \varepsilon$  if

$$T = \frac{1}{1 - \rho} \log \left( \frac{16\sigma_{\max}^*}{(1 - \rho)^2 \varepsilon} \right)$$

- linear convergence
- iteration complexity independent of condition number  $\kappa$
- geometrically decaying stepsizes

## Convergence theory of Muon

---

simplified Muon :  $U_{t+1} = U_t - \eta_t \text{msign}((U_t U_t^\top - M^*)U_t)$

# Convergence theory of Muon

---

simplified Muon :  $U_{t+1} = U_t - \eta_t \text{msign}((U_t U_t^\top - M^*)U_t)$

## Theorem 2 (heavily over-parameterized ( $k \geq d$ ))

- *stepsize*:  $\eta_t = C_\eta \sqrt{\sigma_{\max}^*} \rho^t$  for  $1/2 \leq \rho < 1$ ,  $C_\eta \sim \text{Unif}(1, 2)$
- *initialization*:  $U_0 = \alpha O$  w/ orthonormal  $O$ ,  $0 < \alpha \leq C_\eta \sqrt{\lambda_{\max}^*}$

# Convergence theory of Muon

---

simplified Muon :  $U_{t+1} = U_t - \eta_t \text{msign}((U_t U_t^\top - M^*)U_t)$

## Theorem 2 (heavily over-parameterized ( $k \geq d$ ))

- *stepsize*:  $\eta_t = C_\eta \sqrt{\sigma_{\max}^*} \rho^t$  for  $1/2 \leq \rho < 1$ ,  $C_\eta \sim \text{Unif}(1, 2)$
- *initialization*:  $U_0 = \alpha O$  w/ orthonormal  $O$ ,  $0 < \alpha \leq C_\eta \sqrt{\lambda_{\max}^*}$

With prob. 1, one has  $\|U_T U_T^\top - M^*\| \leq \varepsilon$  if

$$T \geq \frac{1}{1 - \rho} \log \left( \frac{8\sigma_{\max}^*}{\varepsilon} \right)$$

# Convergence theory of Muon

simplified Muon :  $U_{t+1} = U_t - \eta_t \text{msign}((U_t U_t^\top - M^*)U_t)$

## Theorem 2 (heavily over-parameterized ( $k \geq d$ ))

- *stepsize*:  $\eta_t = C_\eta \sqrt{\sigma_{\max}^*} \rho^t$  for  $1/2 \leq \rho < 1$ ,  $C_\eta \sim \text{Unif}(1, 2)$
- *initialization*:  $U_0 = \alpha O$  w/ orthonormal  $O$ ,  $0 < \alpha \leq C_\eta \sqrt{\lambda_{\max}^*}$

With prob. 1, one has  $\|U_T U_T^\top - M^*\| \leq \varepsilon$  if

$$T \geq \frac{1}{1 - \rho} \log \left( \frac{8\sigma_{\max}^*}{\varepsilon} \right)$$

- linear convergence
- iteration complexity independent of condition number  $\kappa$
- geometrically decaying stepsizes

## Comparison with other optimizers?

---

- GD: needs  $\kappa \log \frac{1}{\epsilon}$  iterations

## Comparison with other optimizers?

---

- GD: needs  $\kappa \log \frac{1}{\epsilon}$  iterations
- SignGD (simplified Adam): needs  $\Omega(\kappa)$  iterations

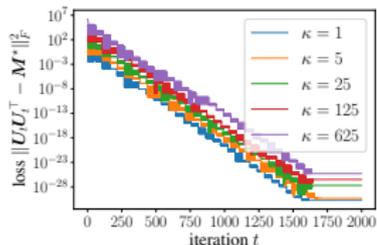
### Theorem 3 (lower bound for SignGD (simplified Adam))

Consider SignGD w/ non-increasing stepsizes and  $\eta_0 \leq 1/16$ . There exists a problem instance such that: for any  $\epsilon \lesssim \kappa^{-2}$ , SignGD w/ warm start needs at least

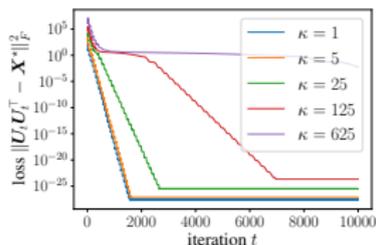
$$\frac{\kappa - 1}{4} \text{ iterations to yield } \epsilon \text{ accuracy}$$

algorithm	parameterization	# iterations	paper
simplified Muon	exactly-parameterized	$\log \frac{1}{\epsilon}$	this work
	over-parameterized	$\log \frac{1}{\epsilon}$	this work
GD	exactly-parameterized	$\kappa \log \frac{1}{\epsilon}$	Chi et al. '19
	over-parameterized	$\kappa^3 \log \frac{1}{\epsilon}$	Stoger et al. '21
	lower bound	$\kappa \log \frac{1}{\epsilon}$	folklore
SignGD	lower bound	$\kappa$	this work

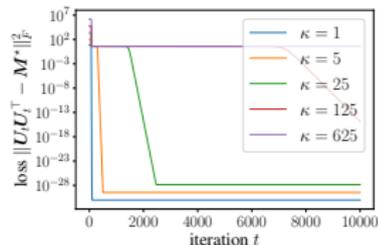
algorithm	parameterization	# iterations	paper
simplified Muon	exactly-parameterized	$\log \frac{1}{\epsilon}$	this work
	over-parameterized	$\log \frac{1}{\epsilon}$	this work
GD	exactly-parameterized	$\kappa \log \frac{1}{\epsilon}$	Chi et al. '19
	over-parameterized	$\kappa^3 \log \frac{1}{\epsilon}$	Stoger et al. '21
	lower bound	$\kappa \log \frac{1}{\epsilon}$	folklore
SignGD	lower bound	$\kappa$	this work



(a) Muon



(b) SignGD



(c) GD

## Analysis: scalar case ( $d = 1$ )

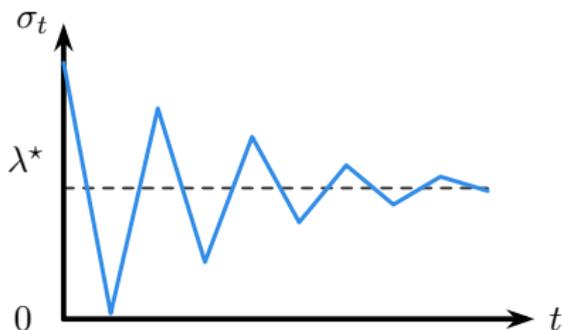
---

$$U_{t+1} = U_t - \eta_t \operatorname{msign}((U_t U_t^\top - M^*)U_t), \quad t = 0, 1, \dots$$

## Analysis: scalar case ( $d = 1$ )

---

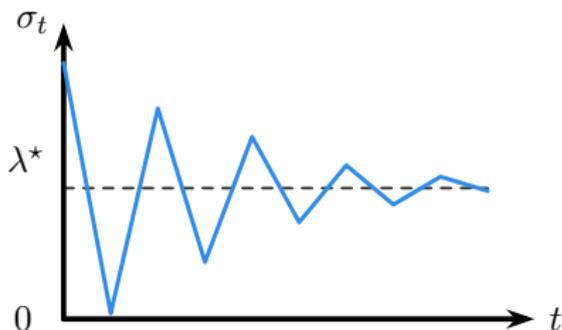
$$\sigma_{t+1} = \sigma_t - \eta_t \text{sign}(\sigma_t^3 - \lambda^* \sigma_t), \quad t = 0, 1, \dots$$



## Analysis: scalar case ( $d = 1$ )

---

$$\sigma_{t+1} = \sigma_t - \eta_t \text{sign}(\sigma_t^3 - \lambda^* \sigma_t), \quad t = 0, 1, \dots$$

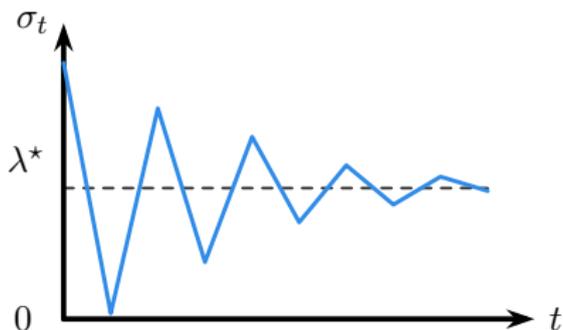


**linear convergence:** using  $\eta_t = C\sqrt{|\lambda^*|}\rho^t$  for  $C \geq 1$ ,  $\rho \in [1/2, 1)$ :

## Analysis: scalar case ( $d = 1$ )

---

$$\sigma_{t+1} = \sigma_t - \eta_t \text{sign}(\sigma_t^3 - \lambda^* \sigma_t), \quad t = 0, 1, \dots$$



**linear convergence:** using  $\eta_t = C\sqrt{|\lambda^*|}\rho^t$  for  $C \geq 1$ ,  $\rho \in [1/2, 1)$ :

$$|\sigma_{t+1}^2 - \lambda^*| \lesssim |\lambda^*|\rho^t \quad (\text{constant conv. rate})$$

## Analysis: the case w/ exact subspace alignment

---

$$M^* = V^* \Lambda^* V^{*\top}$$

## Analysis: the case w/ exact subspace alignment

---

$$M^* = V^* \Lambda^* V^{*\top}$$

if  $U_0 = \underbrace{V^*}_{\text{exact subspace}} \Sigma_0 R^\top$  for ortho  $R$ , then

## Analysis: the case w/ exact subspace alignment

---

$$M^* = \mathbf{V}^* \Lambda^* \mathbf{V}^{*\top}$$

if  $U_0 = \underbrace{\mathbf{V}^* \Sigma_0}_{\text{exact subspace}} \mathbf{R}^\top$  for ortho  $\mathbf{R}$ , then

$$\begin{aligned} U_{t+1} &= U_t - \eta_t \text{msign}((U_t U_t^\top - M^*) U_t) \\ &= \underbrace{\mathbf{V}^* \Sigma_t \mathbf{R}^\top - \eta_t \mathbf{V}^* \text{diag}\{\text{sign}(\Sigma_t^3 - \Lambda^* \Sigma_t)\}}_{\text{subspace is preserved}} \mathbf{R}^\top \end{aligned}$$

# Analysis: the case $w/$ exact subspace alignment

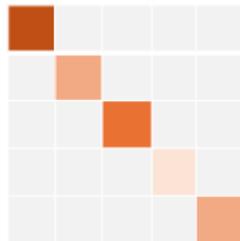
$$M^* = \mathbf{V}^* \Lambda^* \mathbf{V}^{*\top}$$

if  $U_0 = \underbrace{\mathbf{V}^* \Sigma_0}_{\text{exact subspace}} \mathbf{R}^\top$  for ortho  $\mathbf{R}$ , then

$$\begin{aligned} U_{t+1} &= U_t - \eta_t \text{msign}((U_t U_t^\top - M^*) U_t) \\ &= \underbrace{\mathbf{V}^* \Sigma_t \mathbf{R}^\top - \eta_t \mathbf{V}^* \text{diag}\{\text{sign}(\Sigma_t^3 - \Lambda^* \Sigma_t)\}}_{\text{subspace is preserved}} \mathbf{R}^\top \end{aligned}$$

**spectral dynamics:**

$$\Sigma_{t+1} = \Sigma_t - \eta_t \text{diag}\{\text{sign}(\Sigma_t^3 - \Lambda^* \Sigma_t)\}$$



## Analysis: the case w/ exact subspace alignment

---

$$\Sigma_{t+1} = \Sigma_t - \eta_t \text{diag}\{\text{sign}(\Sigma_t^3 - \Lambda^* \Sigma_t)\}$$


$$\sigma_{1,t+1} = \sigma_{1,t} - \eta_t \text{sign}(\sigma_{1,t}^3 - \lambda_1^* \sigma_{1,t})$$

$$\sigma_{2,t+1} = \sigma_{2,t} - \eta_t \text{sign}(\sigma_{2,t}^3 - \lambda_2^* \sigma_{2,t})$$

$$\vdots$$
$$\vdots$$

Spectral dynamics of Muon decouple into  $d$  scalar sequences:

# Analysis: the case $w/$ exact subspace alignment

$$\Sigma_{t+1} = \Sigma_t - \eta_t \text{diag}\{\text{sign}(\Sigma_t^3 - \Lambda^* \Sigma_t)\}$$

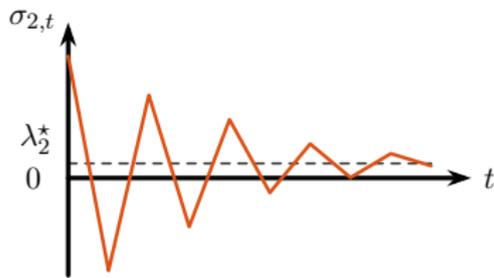
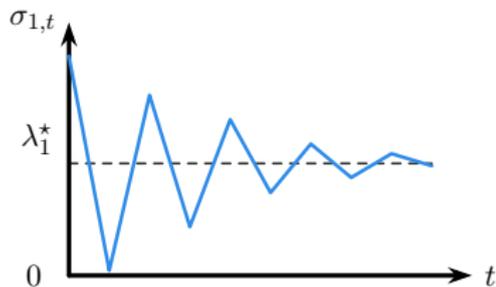
$$\sigma_{1,t+1} = \sigma_{1,t} - \eta_t \text{sign}(\sigma_{1,t}^3 - \lambda_1^* \sigma_{1,t})$$

$$\sigma_{2,t+1} = \sigma_{2,t} - \eta_t \text{sign}(\sigma_{2,t}^3 - \lambda_2^* \sigma_{2,t})$$

$$\vdots$$
$$\vdots$$

Spectral dynamics of Muon decouple into  $d$  scalar sequences:

- each sequence converges linearly w/ constant conv. rates



## Analysis: more general case

---

Subspace alignment?

## Analysis: more general case

---

Subspace alignment?

- always guaranteed when heavily over-parameterized ( $k \geq d$ )

## Analysis: more general case

---

Subspace alignment?

- always guaranteed when heavily over-parameterized ( $k \geq d$ )
- approximately satisfied under small initialization ( $r \leq k < d$ )

$$U_1 \approx \cancel{U_0} - \eta_0 \text{msign}(\cancel{U_0 U_0^\top U_0} - \mathbf{M}^* U_0)$$

## Analysis: more general case

---

Subspace alignment?

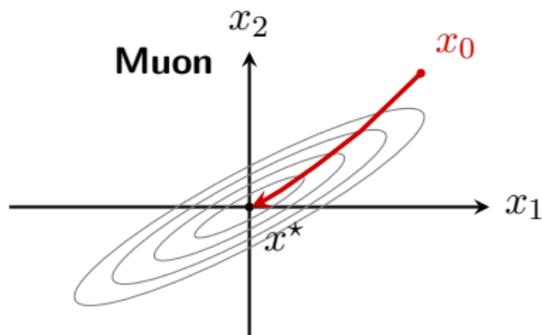
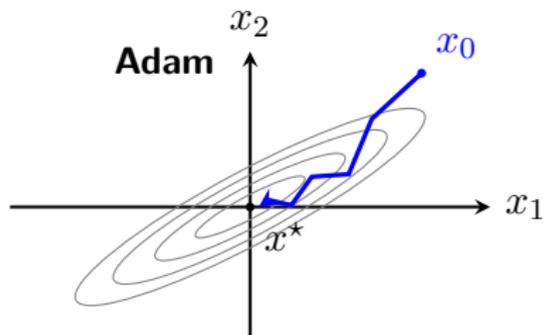
- always guaranteed when heavily over-parameterized ( $k \geq d$ )
- approximately satisfied under small initialization ( $r \leq k < d$ )

$$U_1 \approx \cancel{U_0} - \eta_0 \text{msign}(\cancel{U_0 U_0^\top} U_0 - M^* U_0)$$

- **challenge:** avoid blow-up of approx. error over time

## Concluding remarks

---



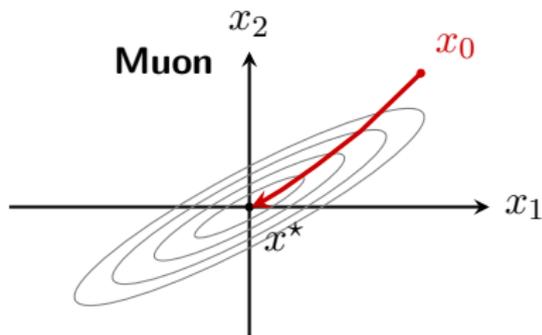
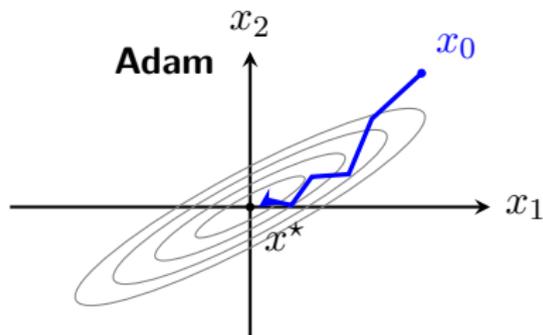
For matrix factorization & ICL of linear transformers:

- **Muon**: condition-number-free fast convergence
  - spectral dynamics: decouple into independent scalar sequences
- **Adam**: slow convergence if Hessian is not nearly diagonal

“Preconditioning benefits of spectral orthogonalization in Muon,” J. Ma, Y. Huang, Y. Chi, Y. Chen, arXiv:2601.13474, 2026

# Concluding remarks

---



## Future directions:

- understand benefits of momentum
- more case studies (recently, Gonon et al. '26, Li et al. '26, ...)
- general theory that demonstrates benefits of Muon

“Preconditioning benefits of spectral orthogonalization in Muon,” J. Ma, Y. Huang, Y. Chi, Y. Chen, arXiv:2601.13474, 2026