

# Exploiting Weight-Space Symmetries for Approximating Curvature

Artem Artemev, Rui Xia, Benjamin M. Boyd, Youjing Yu, Felix Dangel,  
Guillaume Hennequin, Alberto Bernacchia

June 2026

# Curvature is useful but expensive

## Curvature matters for:

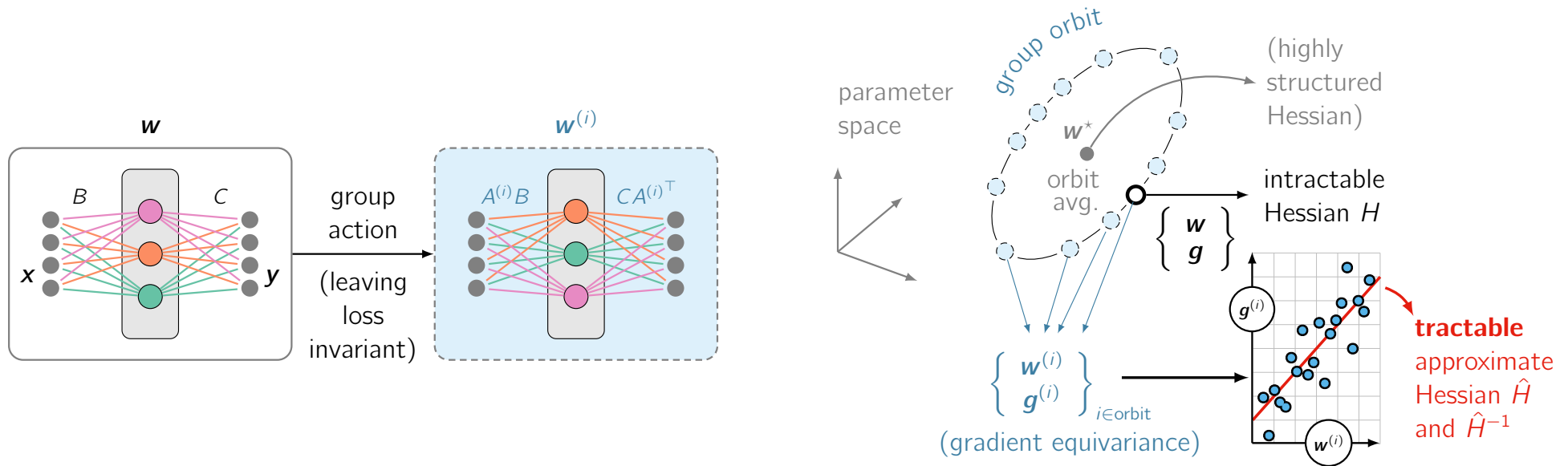
- Second-order optimization
- Bayesian inference
- Continual learning
- Pruning / compression

**The problem:** Full Hessian is  $O(D^2)$

*Our focus is on second-order optimization where KFAC, Shampoo/Muon work well.*

*We introduce a framework based on weight-space symmetries.*

# Loss invariance $\rightarrow$ gradient equivariance $\rightarrow$ curvature



*A single gradient reveals curvature along the entire orbit.*

## Method: orbit-averaged Hessian

**Taylor expand** around orbit average  $w^*$ :

$$g - g^* \approx H^*(w - w^*)$$

**Average** the secant equation over the orbit:

$$S_g \approx H^* S_w H^*$$

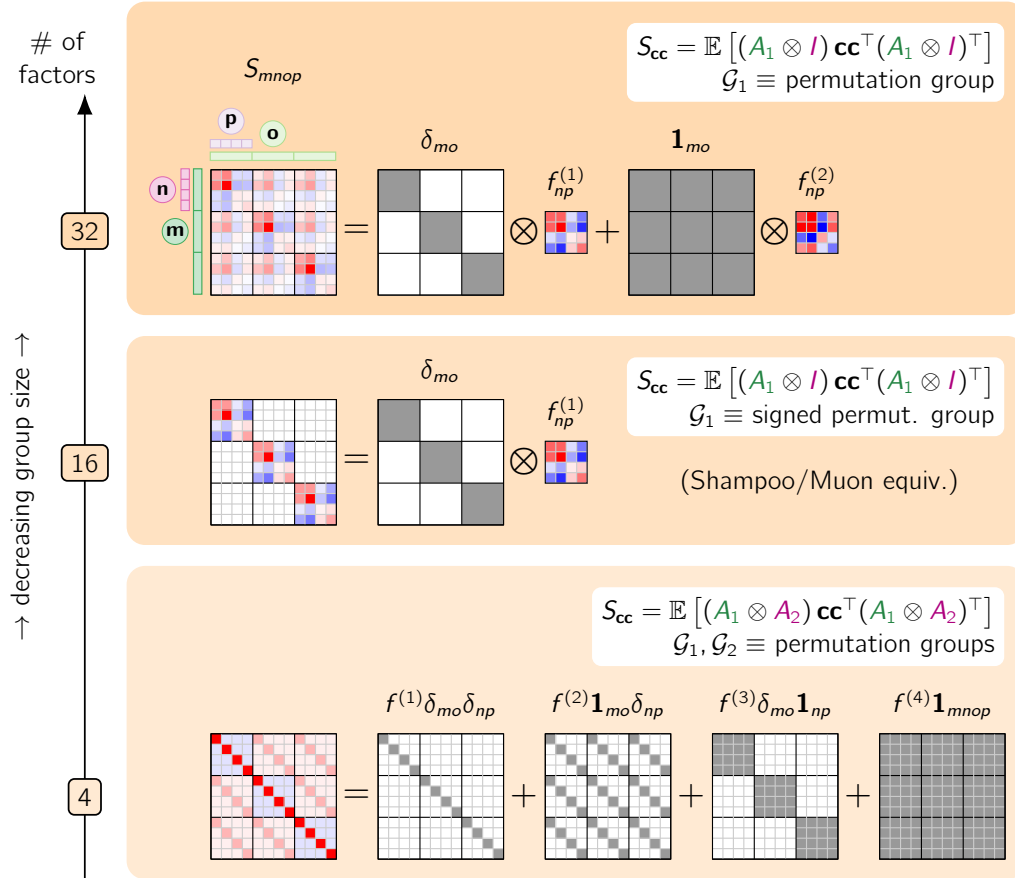
$\left\{ \begin{array}{l} \mathcal{G} \text{ is large} \rightarrow \text{cheap} \\ \mathcal{G} \text{ is small} \rightarrow \text{accurate} \end{array} \right.$

**Solve** for the positive semi-definite curvature:

$$H_{\text{PD}}^* = S_w^{-\frac{1}{2}} \left( S_w^{\frac{1}{2}} S_g S_w^{\frac{1}{2}} \right)^{\frac{1}{2}} S_w^{-\frac{1}{2}}$$

$$\hat{H}_g = S_g^{1/2}$$

# Orbit average structure depends on group size



Larger group  $\rightarrow$  fewer basis elements  $\rightarrow$  cheaper.  
 Shampoo/Muon sits in the middle.

## Symo recovers Shampoo/Muon as a Special case

When transformation groups alternate with identity maps, e.g. for an MLP:

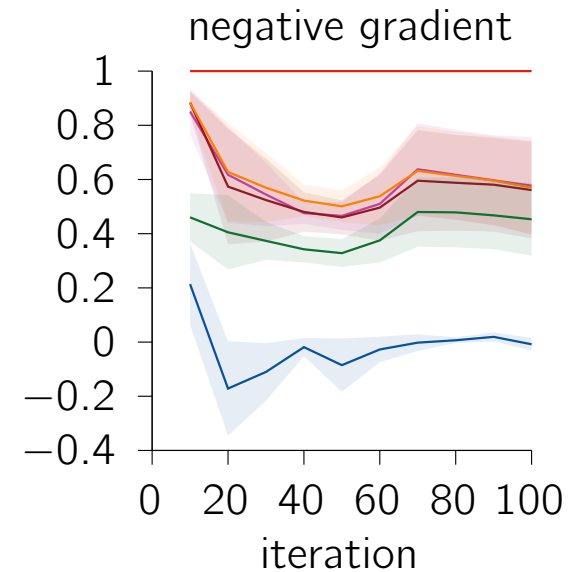
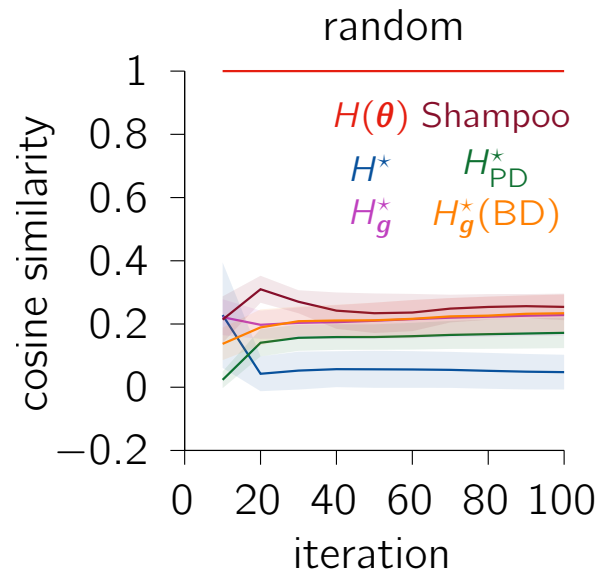
$$\mathcal{G} \stackrel{\text{def}}{=} I \times \mathcal{B}_1 \times I \times \mathcal{B}_2 \times \dots$$

the Symo update coincides with Muon/Shampoo:

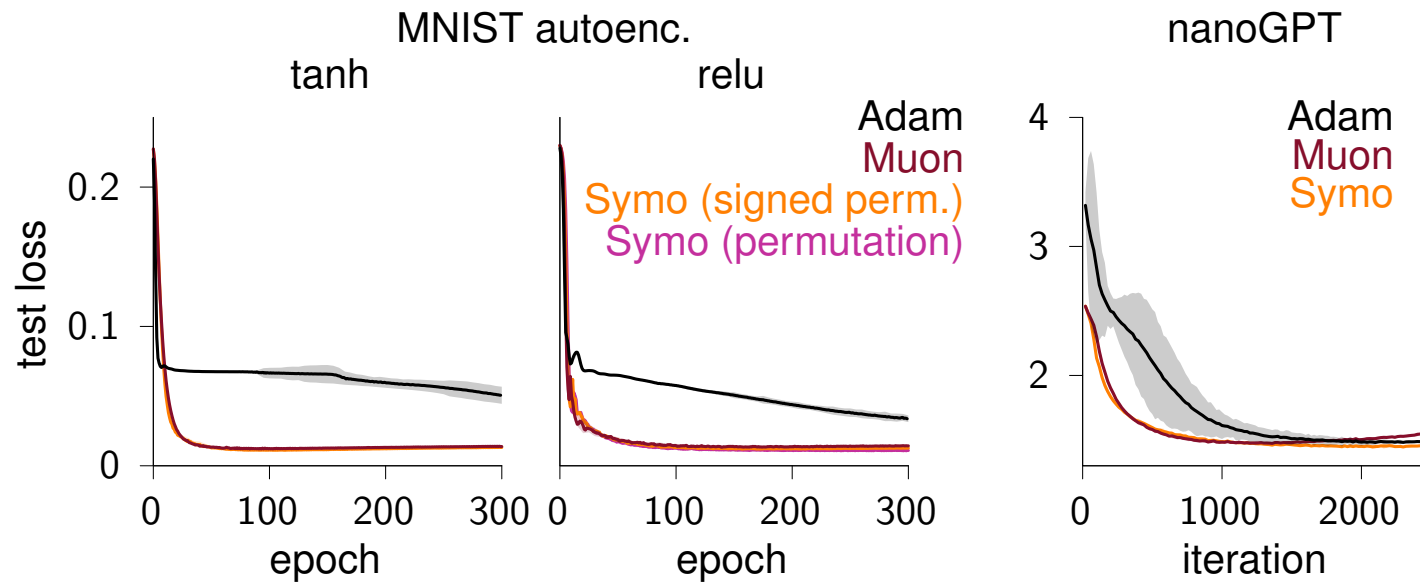
$$\Delta G = \sqrt{d} G (G^\top G)^{-1/2}$$

# Hessian approximation quality

- Along *gradient direction*,  $\hat{H}_g$  achieves 0.4 – 0.8 cosine similarity
- $\hat{H}_g^{\text{BD}}$  overlaps with Shampoo ✓
- $\hat{H}_g$  beats  $\hat{H}_{\text{PD}}$  (rank deficiency of  $S_w$ )

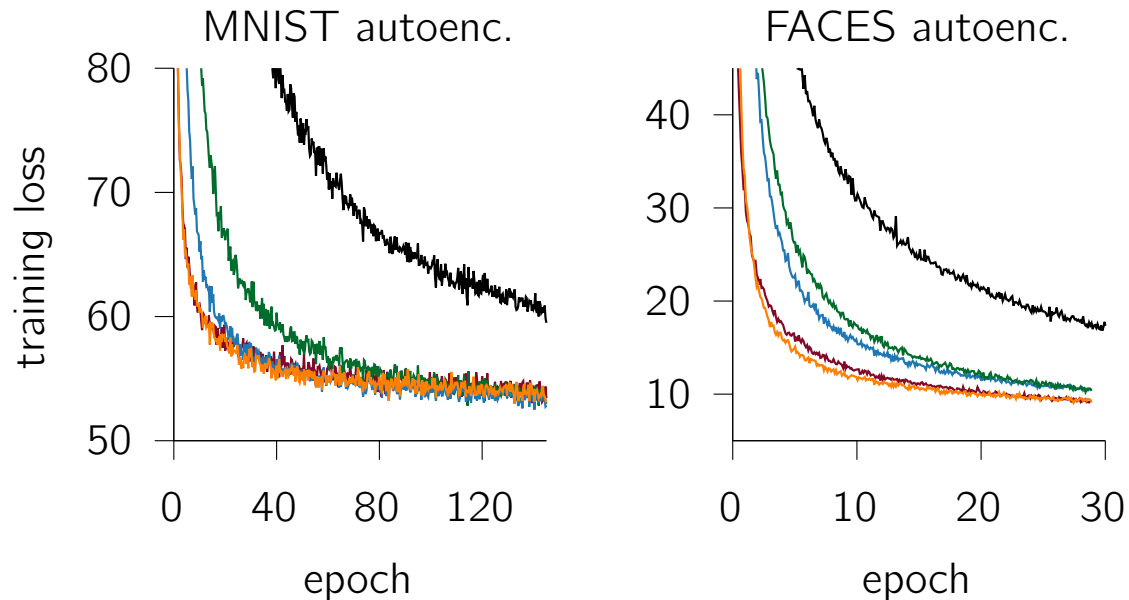


# Empirical equivalence



# Exploring different symmetry groups

Adam   Shampoo   Symo  
Symo (larger group)   Symo (smaller group, block-diag)



- *Larger group*: outperforms Adam (only 64 free parameters!)
- *Smaller group, block-diag*: matches Shampoo
- *Symo*: between the two

The group choice systematically controls the accuracy–cost trade-off.

# Summary

## Contributions

- Hessian approximation from *a single gradient* via orbit averaging.
- Complexity controlled by *group size*.
- Shampoo and Muon recovered as *special cases*.
- Validated on MLPs and Transformers

## Future directions

- Bayesian inference, continual learning, pruning
- Engineering:
  - Automated model symmetry-to-optimizer compiler
  - Scaling to larger models

[github.com/mtkresearch/symm\\_opt](https://github.com/mtkresearch/symm_opt)