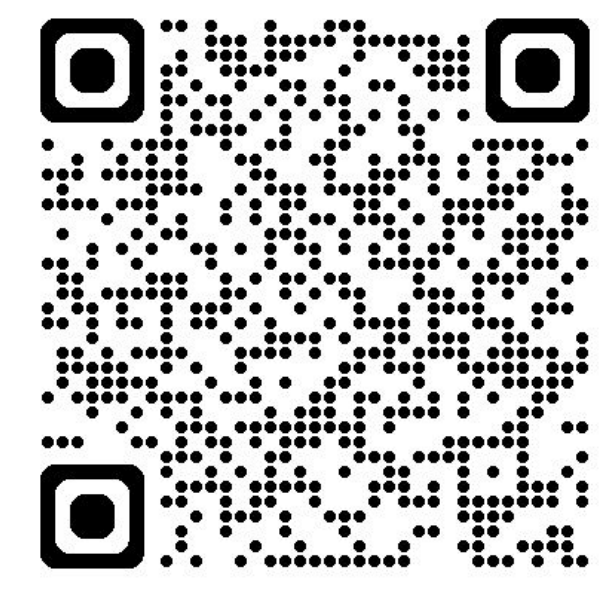


# Adam Reduces a Unique Form of Sharpness: Theoretical Insights Near the Minimizer Manifold

Xinghan Li\* Haodong Wen\* Kaifeng Lyu†

Institute for Interdisciplinary Information Sciences, Tsinghua University

\*Equal contribution; alphabet ordering †Corresponding author



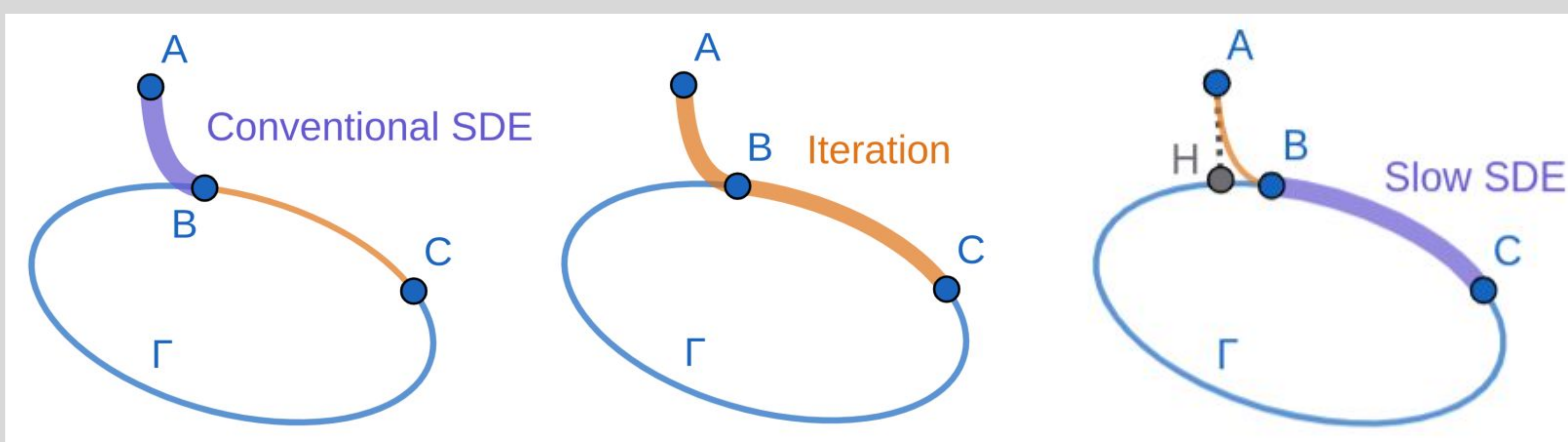
## Key Contributions

- We capture the behaviour of Adam near the minimizer manifold for **up to  $O(\eta^{-2})$  time** using SDE.
- Our SDE shows that Adam implicitly minimizes **a new kind of sharpness:  $\text{tr}(\text{diag}(\mathbf{H})^{1/2})$**  with label noise, instead of SGD's familiar  $\text{tr}(\mathbf{H})$ . To highlight this discrepancy, we conducted two case studies: A sparse linear regression task with diagonal linear networks where Adam has provable generalization benefit over SGD, and one matrix factorization task where Adam generalizes worse.

## Motivation

- SGD's implicit bias toward flatter minima is fairly well understood, but modern-day training uses Adam more.
- We lack a rigorous account of **what kind of sharpness Adam really pursues** and whether that helps or hurts generalization.

## Method: Slow SDE



The Slow SDE **separates the slow implicit bias motion from the fast convergence motion**, and capture the slow motion (moving at speed  $\eta^2$ ).

## Main Results

A Class of Adaptive Gradient Methods

$$\begin{aligned} \mathbf{m}_{k+1} &:= \beta_1 \mathbf{m}_k + (1 - \beta_1) \nabla \ell_k(\boldsymbol{\theta}_k) \\ \mathbf{v}_{k+1} &:= \beta_2 \mathbf{v}_k + (1 - \beta_2) V(\nabla \ell_k(\boldsymbol{\theta}_k) \nabla \ell_k(\boldsymbol{\theta}_k)^\top) \\ \boldsymbol{\theta}_{k+1} &:= \boldsymbol{\theta}_k - \eta S(\mathbf{v}_{k+1}) \mathbf{m}_{k+1}. \end{aligned}$$

Table 1: Examples of  $V, S$  functions for some optimizers in the AGM Framework.

| Optimizer        | Function $V$  | Function $S$   | Remarks   |
|------------------|---|--|---|
| Adam             | $V(\mathbf{M}) = \text{diag}(\mathbf{M})$                     | $S(\mathbf{v}) = \text{Diag}(1/(\sqrt{\mathbf{v}} + \epsilon))$        |   |
| Adam-mini        | $V(\mathbf{M})_i = \frac{1}{ B(i) } \sum_{j \in B(i)} M_{jj}$ | $S(\mathbf{v}) = \text{Diag}(1/(\sqrt{\mathbf{v}} + \epsilon))$        | Parameters partitioned; $i$ belongs to block $B(i)$ . $i$ belongs to layer $L(i)$ in the model. |
| Adalayer         | $V(\mathbf{M})_i = \frac{1}{ L(i) } \sum_{j \in L(i)} M_{jj}$ | $S(\mathbf{v}) = \text{Diag}(1/(\sqrt{\mathbf{v}} + \epsilon))$        |   |
| AdamE- $\lambda$ | $V(\mathbf{M}) = \text{diag}(\mathbf{M})$                     | $S(\mathbf{v}) = \text{Diag}(1/(\mathbf{v}^\odot \lambda + \epsilon))$ |   |

## Slow SDE for AGMs:

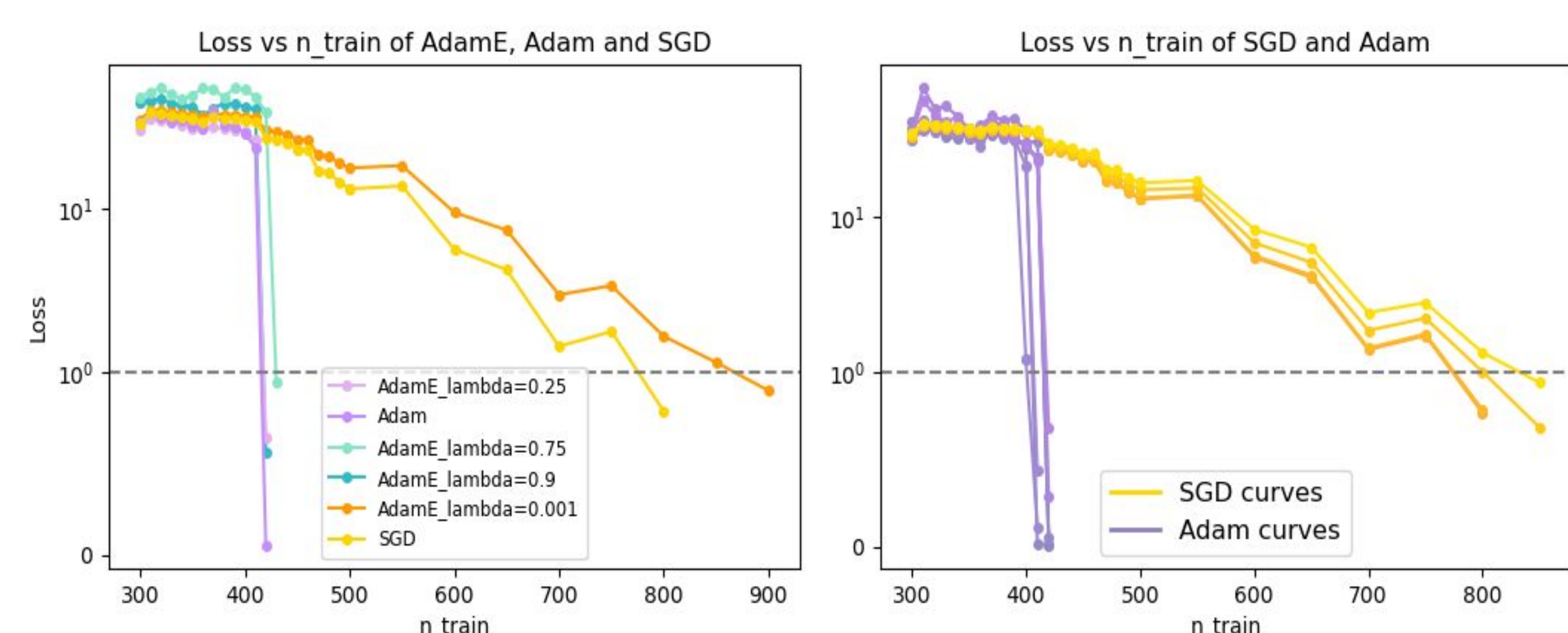
$$\begin{cases} d\zeta(t) = P_{\zeta, \mathbf{S}(t)}(\underbrace{\Sigma^{1/2}(\zeta(t); \mathbf{S}(t)) d\mathbf{W}_t}_{\text{diffusion}} - \underbrace{\frac{1}{2} \mathbf{S}(t) \nabla^3 \mathcal{L}(\zeta) [\Sigma_\diamond(\zeta(t); \mathbf{S}(t))] dt}_{\text{drift}}), \\ d\mathbf{v}(t) = \underbrace{c(V(\Sigma(\zeta)) - \mathbf{v}) dt}_{\text{Preconditioner drift}}. \end{cases} \quad \left( \frac{1 - \beta_2}{\eta^2} = c, \quad \zeta_0 \in \Gamma \right)$$

- $\Sigma(\zeta)$  is the covariance at  $\zeta$ ,  $\mathbf{S}(t) := \mathbf{S}(\mathbf{v}(t))$ ,  $P_{\zeta, \mathbf{S}(t)}$ :  $\zeta$ 's projection on manifold, defined by gradient flow pre-conditioned with  $\mathbf{S}(t)$ .  $\Sigma_\parallel(\zeta; \mathbf{S})$  and  $\Sigma_\diamond(\zeta; \mathbf{S})$  are matrices related to  $\Sigma(\zeta)$ ,  $\zeta$  and  $\mathbf{S}$ .
- The drift term in Slow SDE can be interpreted as **adaptive semi-gradient descent** minimizing  $\mu(\zeta, \mathbf{v}) := \langle \nabla^2 \mathcal{L}(\zeta), \Sigma_\diamond(\zeta(t); \mathbf{S}(t)) \rangle$

## Adam's Generalization Benefit with Label Noise

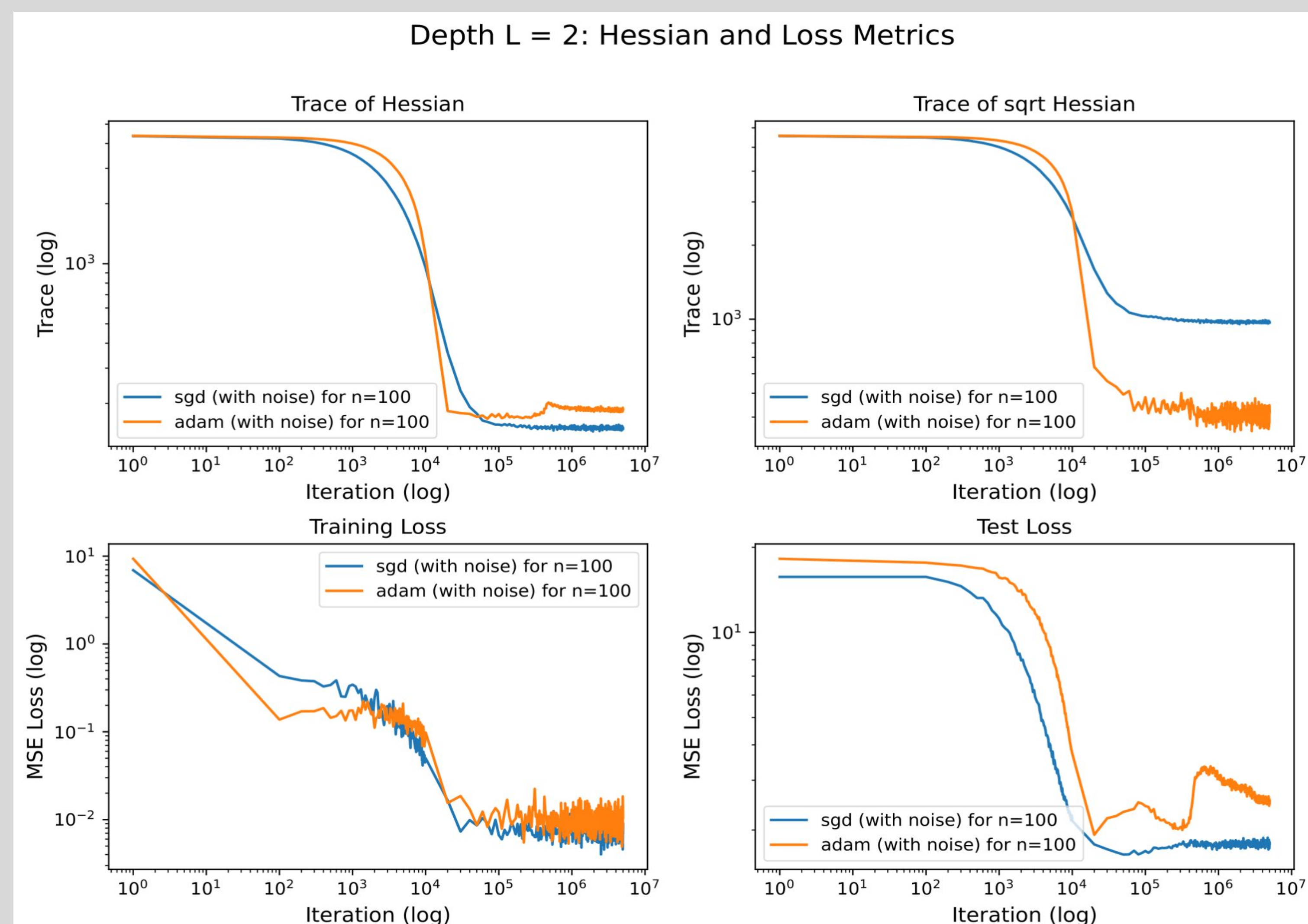
- Label noise  $\Sigma \equiv \alpha \nabla^2 \mathcal{L}$ , SGD ODE:  $d\zeta(t) = -\frac{\alpha}{4} P_\parallel(\zeta) \nabla^3 \mathcal{L}(\zeta) [\mathbf{I}] dt$ .
- **AGM's slow SDE reduces to an ODE:**  $\begin{cases} d\mathbf{v}(t) = c(V(\Sigma(\zeta)) - \mathbf{v}) dt, \\ d\zeta(t) = -\frac{\alpha}{2} S(\mathbf{v}) P_{\parallel, S(\mathbf{v})}(\zeta) S(\mathbf{v}) \nabla^3 \mathcal{L}(\zeta) [S(\mathbf{v})] dt. \end{cases}$
- The fixed point of this ODE must satisfy  $\nabla \text{tr}(\text{Diag}(\mathbf{H})^{1/2}) = 0$
- Changing Adam's sqrt to  $^\wedge \lambda$  results in  $\nabla \text{tr}(\text{Diag}(\mathbf{H})^{1-\lambda}) = 0$

## Experiment: Sparse Regression with Diagonal Net



**Takeaway: Adam's unique implicit bias aligns better with the sparsity requirement in this case** (Adam- $\lambda$  finds the optimum with minimal  $\lambda$ -norm), which arises from taking into consideration the 2nd order momentum compared to SGD.

## Adam Loses in Matrix Factorization



**Takeaway: Adam's implicit bias hurts generalizability in this case.** Adam's implicit regularization differs qualitatively from SGD's, but the benignity of this difference depends.