# WHAT HAPPENS AFTER SGD REACHES ZERO LOSS? – A MATHEMATICAL FRAMEWORK

Zhiyuan Li[†]    Tianhao Wang[‡]    Sanjeev Arora[†]

[†]Princeton University    [‡]Yale University

## INTRODUCTION

Stochastic gradient descent (SGD) is widely used in training of modern machine learning models such as deep neural networks, and the implicit bias of SGD underlies the generalization ability of the trained models. While it still remains unclear how to mathematically characterize such bias.
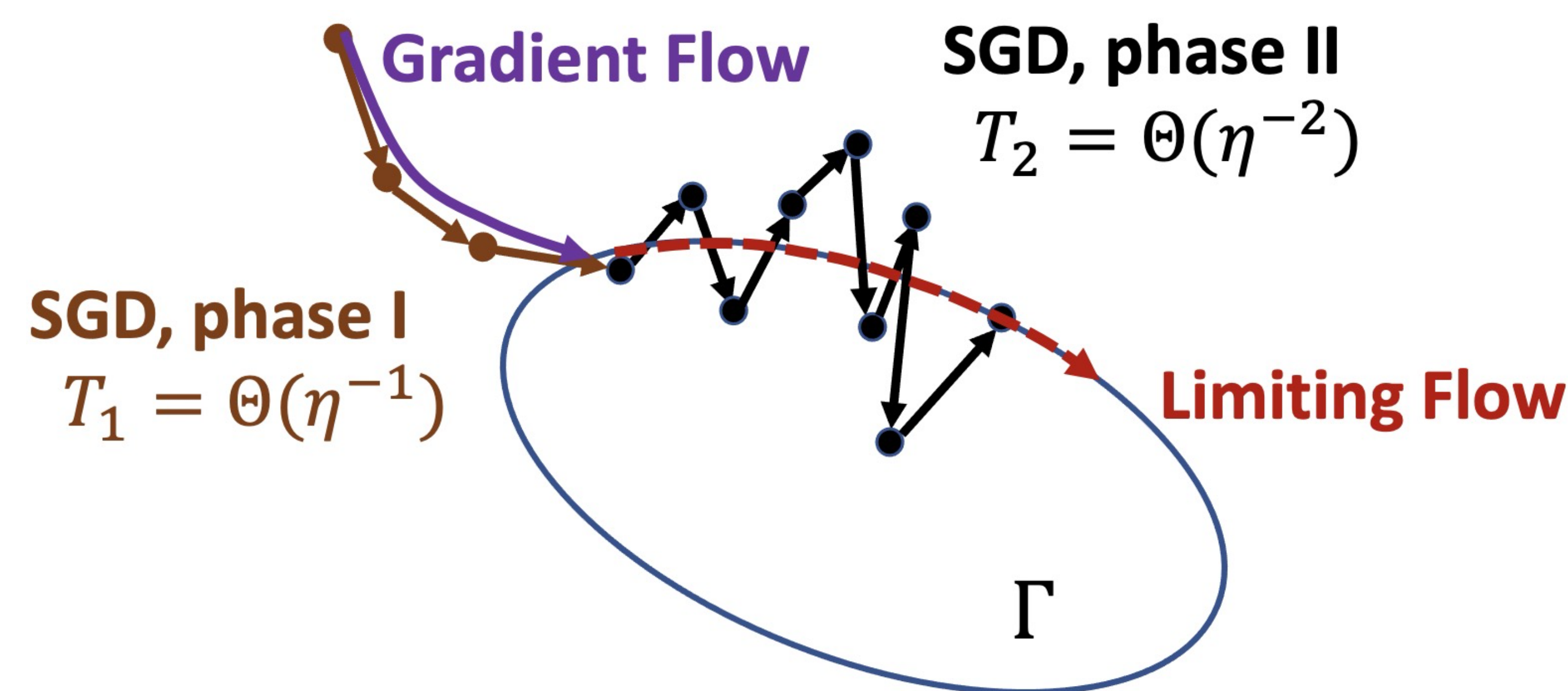
**Formulation of SGD:** Given training loss $L : \mathbb{R}^D \to \mathbb{R}$,

$$x_\eta(k+1) = x_\eta(k) - \eta(\nabla L(x_\eta(k)) + \sqrt{\Xi}\sigma_{\xi_k}(x_\eta(k))) \quad (1)$$

- $\eta$ is the learning rate (LR)
- $\sigma(x) = [\sigma_1(x), \sigma_2(x), \ldots, \sigma_\Xi(x)] \in \mathbb{R}^{D \times \Xi}$ is the noise function
- $\xi_k$ is sampled uniformly from $\{1, 2, \ldots, \Xi\}$ and $\mathbb{E}_{\xi_k}[\sigma_{\xi_k}(x)] = 0$

**Main Contributions of This Work:**

#1. A mathematical framework to study implicit bias of SGD with small LR

#2. Provable generalization benefit of stochasticity: Minimax optimal rate for learning sparse quadratically overparametrized linear models.



## PROBLEM SETTING

**Manifold of Local Minimizers:** $\Gamma$ is a $(D-M)$-dimensional submanifold of $\mathbb{R}^D$ such that for all $x \in \Gamma$, $x$ is a local minimizer of $L$ and $\text{rank}(\nabla^2 L(x)) = M$.

**When Does Such A Manifold Exist?** Overparametrization!

**Canonical SDE Approximation of SGD:**

$$d\widetilde{X}_\eta(t) = -\eta \nabla L(\widetilde{X}_\eta(t))dt + \eta \cdot \sigma(\widetilde{X}_\eta(t))dW(t) \quad (2)$$

where $\{W(t)\}_{t\geq 0}$ is a $\Xi$-dimensional Wiener Process and $\Sigma(x) = \sigma(x)\sigma(x)^\top$ is the covariance matrix of gradient noise.
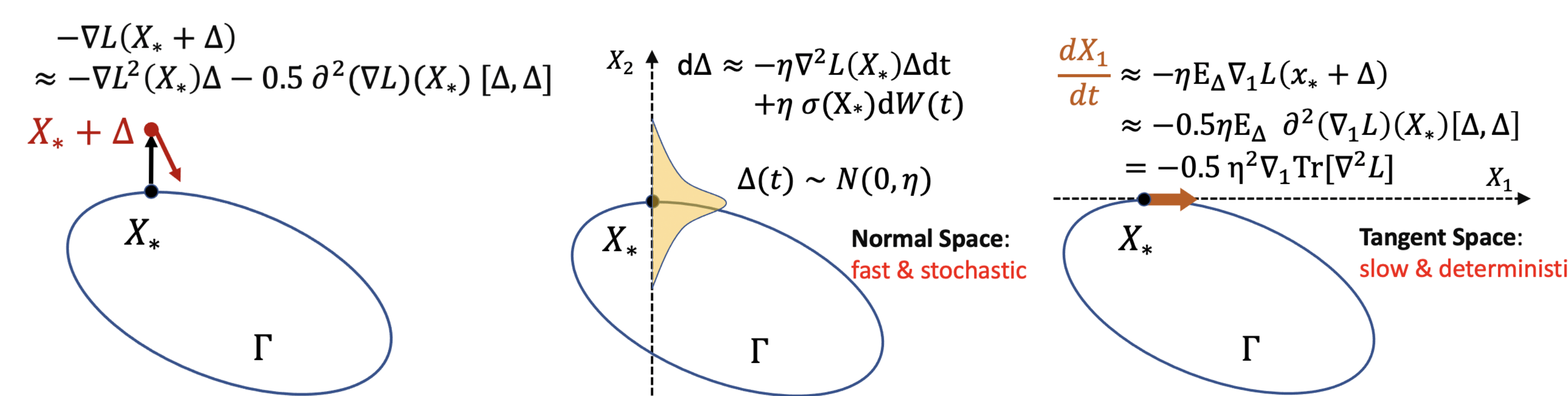
## INTUITIVE EXPLANATION OF THE IMPLICIT BIAS

Blanc et. al, (2020) showed that around some manifold of local minimizers SGD decreases $\text{tr}[\nabla^2 L]$ if gradient covariance is equal to trace of hessian, $\Sigma \equiv \nabla^2 L$, on the manifold. (e.g., SGD with label noise)

**Taylor Expansion Around A Local Minimizer:** Let $\Delta(t) = \widetilde{X}_\eta(t) - X^*$,

$$d\Delta(t) \approx -\eta \nabla^2 L(X^*)\Delta(t)dt + \eta\sigma(X^*)dW(t)$$

which behaves like an Ornstein-Uhlenbeck process in the normal space



The fast dynamics in the normal space activates the second order Taylor expansion in the tangent space, creating a $\Theta(\eta^2)$ velocity, which is slow but deterministic and acumulates over time.

**Issue:** The above local analysis only holds for $O(\eta^{-1.6})$ time. How to do global analysis?

## OUR APPROACH: SEPARATING SLOW FROM FAST

**Time Rescaling for SDE:** Let $X_\eta(t) = \widetilde{X}_\eta(t/\eta^2)$, then

$$dX_\eta(t) = \underbrace{-\eta^{-1}\nabla L(X_\eta(t))dt}_{\text{Fast}} + \underbrace{\sigma(X_\eta(t))dW(t)}_{\text{Slow}}$$

As $\eta \to 0$, the *Fast* part rapidly drives $X_\eta(t)$ towards $\Gamma$ via the projection induced by the gradient flow, denoted by $\Phi(X_\eta(t))$.

**Lemma 1.** $\partial\Phi(x)\nabla L(x) \equiv 0$.

Applying Lemma 1 and Ito's lemma, we get

$$d\Phi(X_\eta(t)) = \partial\Phi(X_\eta)\sigma(X_\eta)dW(t) + \frac{1}{2}\partial^2\Phi(X_\eta)[\sigma(X_\eta)\sigma(X_\eta)^\top]dt.$$

Since $\Phi(X_\eta(t)) \approx X_\eta(t)$ near $\Gamma$, the *Slow* part survives:

$$dX_\eta(t) \approx \underbrace{\partial\Phi(X_\eta)\sigma(X_\eta)dW(t)}_{\text{Tangent noise}} + \underbrace{\frac{1}{2}\partial^2\Phi(X_\eta)[\sigma(X_\eta)\sigma(X_\eta)^\top]dt}_{\text{Compensation and regularization}}.$$

The above analysis can be made rigorous and extended to SGD by viewing SGD as an asymptotically continuous stochastic process and further applying the classic results by Katzenberger (1991).

## MAIN RESULTS

**Lemma 2.** $\partial\Phi(x)$ is the projection matrix of the tangent space of $\Gamma$ at $x$.

**Notation:** . Define $\Sigma_\parallel(x) = \partial\Phi(x)\Sigma(x)\Phi(x)$ (noise covariance in the tangent space), $\Sigma_\perp(x) = (I - \partial\Phi(x))\Sigma(x)(I - \partial\Phi(x))$ (noise covariance in the normal space), and $\Sigma_{\parallel,\perp} = \Sigma_{\perp,\parallel}^\top = \partial\Phi(x)\Sigma(x)(I - \partial\Phi(x))$ (covariance across the tangent and normal space).

**Lyapunov Operator:** For a symmetric matrix $H$, define $W_H = \{\Sigma \mid \Sigma = \Sigma^\top, HH^\dagger\Sigma = \Sigma = \Sigma HH^\dagger\}$. The *Lyapunov operator* $\mathcal{L}_H : W_H \to W_H$ is defined as $\mathcal{L}_H(\Sigma) = H\Sigma + \Sigma H$.

**Main Theorem.** For SGD (1) and any $T > 0$, $x_\eta(\lfloor T/\eta^2 \rfloor)$ converges in distribution to $Y(T)$ as $\eta \to 0$, where $Y(T)$ is the solution to the following SDE at time $T$ when the global solution exists:

$$dY(t) = \underbrace{\Sigma_\parallel^{1/2}(Y)dW(t)}_{\text{Tangent noise}} - \underbrace{\frac{1}{2}\nabla^2 L(Y)^\dagger\partial^2(\nabla L)(Y)[\Sigma_\parallel(Y)]dt}_{\text{Tangent noise compensation}}$$

$$- \underbrace{\frac{1}{2}\partial\Phi(Y)\partial^2(\nabla L)(Y)[\nabla^2 L(Y)^\dagger\Sigma_{\perp,\parallel}(Y)]dt}_{\text{Mixed regularization}}$$

$$- \underbrace{\frac{1}{2}\partial\Phi(Y)\partial^2(\nabla L)(Y)[\mathcal{L}_{\nabla^2 L}^{-1}(\Sigma_\perp(Y))]dt}_{\text{Normal regularization}}.$$

## PROVABLE GENERALIZATION BENEFIT OF STOCHASTICITY

**Setting:** Data $\{(z_i, y_i)\}_{i=1}^n$ where $z_1, \ldots, z_n \overset{i.i.d}{\sim} \text{Unif}(\{\pm 1\}^d)$ or $\mathcal{N}(0, I_d)$ and each $y_i = \langle z_i, w^* \rangle$ for some unknown $\kappa$-sparse $w^* \in \mathbb{R}^d$. Denote $x = \binom{u}{v} \in \mathbb{R}^D = \mathbb{R}^{2d}$. For each $i \in [n]$, define $f_i(x) = \langle z_i, u^{\odot 2} - v^{\odot 2} \rangle$. Consider the $\ell_2$ loss $L(x) = \frac{1}{n}\sum_{i=1}^n (f_i(x) - y_i)^2$.

**Label Noise SGD:** At iteration $k$, replace the true label $y_{i_k}$ by a perturbed label $y_{i_k} + \delta_k$ where $\delta_k \sim \text{Unif}(\{\pm 1\})$ and run SGD on the perturbed label.

**Regularizer:** $R(x) = \text{tr}[\nabla^2 L(x)] = \frac{4}{n}\sum_{j=1}^D \left(\sum_{i=1}^n z_{i,j}^2\right)(u_j^2 + v_j^2)$.

**Limiting Dynamics** = Riemannian gradient flow of $R$ on $\Gamma$:

$$dx(t) = -\partial\Phi(x(t))\nabla R(x(t))dt.$$

**Optimal Sparse Recovery** $\Longleftarrow$ Constrained minimization of $R$ on $\Gamma$.

**Theorem.** Under the above setting with $n \geq \Omega(\kappa \ln d)$ data, for any generic initialization $x_0$ and any $\epsilon > 0$, there exist $\eta_0, T > 0$ such that for any $\eta < \eta_0$, label noise SGD with LR $\eta$ returns an $\epsilon$-optimal solution in $\lfloor T/\eta^2 \rfloor$ steps with high probability.