

Anomalous diffusion dynamics of learning in deep neural networks

Guozhang Chen,¹ Cheng Kevin Qu,¹ and Pulin Gong^{1,*}

¹*School of Physics, University of Sydney, Sydney, NSW 2006, Australia*

(Dated: July 27, 2021)

Learning in deep neural networks (DNNs) is implemented through minimizing a highly non-convex loss function, typically by a stochastic gradient descent (SGD) method. This learning process can effectively find good wide minima without being trapped in poor local ones. We present a novel account of how such effective deep learning emerges through the interactions of the SGD and the geometrical structure of the loss landscape. We find that the SGD exhibits rich, complex dynamics when navigating through the loss landscape; initially, the SGD exhibits anomalous superdiffusion, which attenuates gradually and changes to subdiffusion at long times when approaching a solution. Such learning dynamics happen ubiquitously in different DNNs types such as ResNet and VGG-like networks and are insensitive to batch size and learning rate. The anomalous superdiffusion process during the initial learning phase indicates that the motion of SGD along the loss landscape possesses intermittent, big jumps; this non-equilibrium property enables the SGD to escape from sharp local minima. By adapting the methods developed for studying energy landscapes in complex physical systems, we find that such superdiffusive learning dynamics are due to the interactions of the SGD and the fractal-like regions of the loss landscape. We further develop a simple model to demonstrate the mechanistic role of the fractal-like loss landscape in enabling the SGD to effectively find global minima. Our results reveal the effectiveness of deep learning from a novel perspective and have implications for designing efficient deep neural networks.

I. INTRODUCTION

Deep neural networks (DNNs) trained by stochastic gradient descent (SGD) have achieved great success in many application areas [1]. As often assumed, the SGD optimizer of highly non-convex loss functions is rarely trapped in local minima, and effectively finds wide ones with good generalization [2, 3]. Understanding how this property emerges from the DNNs is of fundamental importance for deciphering the secret of the remarkable effectiveness of deep learning [4].

Recently, progress has been made in either characterizing the structure of loss functions or the dynamics of SGD for gaining comprehension of deep learning. For instance, loss functions have been studied by using random matrix theory [5], algebraic geometry methods [6] and visualization-based methods [7]. The dynamics of SGD have been examined via models of stochastic gradient Langevin dynamics with an assumption that gradient noise is Gaussian [8, 9]; in these models, the SGD is assumed to be driven by Brownian motion in particular. However, it has been increasingly realized that such Brownian motion-based characterizations of SGD dynamics are inappropriate, as SGD dynamics commonly exhibit highly anisotropic and dynamic-changing properties [10–13], suggesting the presence of rich, complex learning dynamics in DNNs. Recently, it has been shown that an inverse relation holds between the landscape flatness and the weight variance [14]. Despite the progress made by these studies, the fundamental questions of how the interaction of SGD with the structure of the loss

function gives rise to complex learning dynamics, and whether and how such dynamics enable SGD to find wide minima remain unknown.

In this study, by adapting the methods developed in nonequilibrium physical systems, we find that the interactions of the loss landscape and the SGD give rise to complex learning dynamics; these include anomalous superdiffusion during the initial learning phase, which changes to subdiffusion at long times when approaching a solution. During this process, the SGD walker moves from rougher (fractal-like) regions to flatter regions of the loss landscape. The fractal-like regions of the loss landscape indicate that they possess varying steepness (Fig. 1) and that the corresponding loss gradient displays large fluctuations with heavy-tailed distributions, thus resulting in superdiffusive learning dynamics. Superdiffusion consists of small movements that are intermittently interrupted by big jumps; these enable SGD to escape from local minima, thus effectively exploring the loss landscape. This computational advantage of superdiffusion is further illustrated in a simple model where the SGD interacts with a low-dimensional fractal loss landscape. Due to its movement even slower than a normal diffusive process (i.e. Brownian motion), the subdiffusive dynamics, on the other hand, may consolidate the residence of the SGD in the flatter areas with good solutions.

* pulin.gong@sydney.edu.au

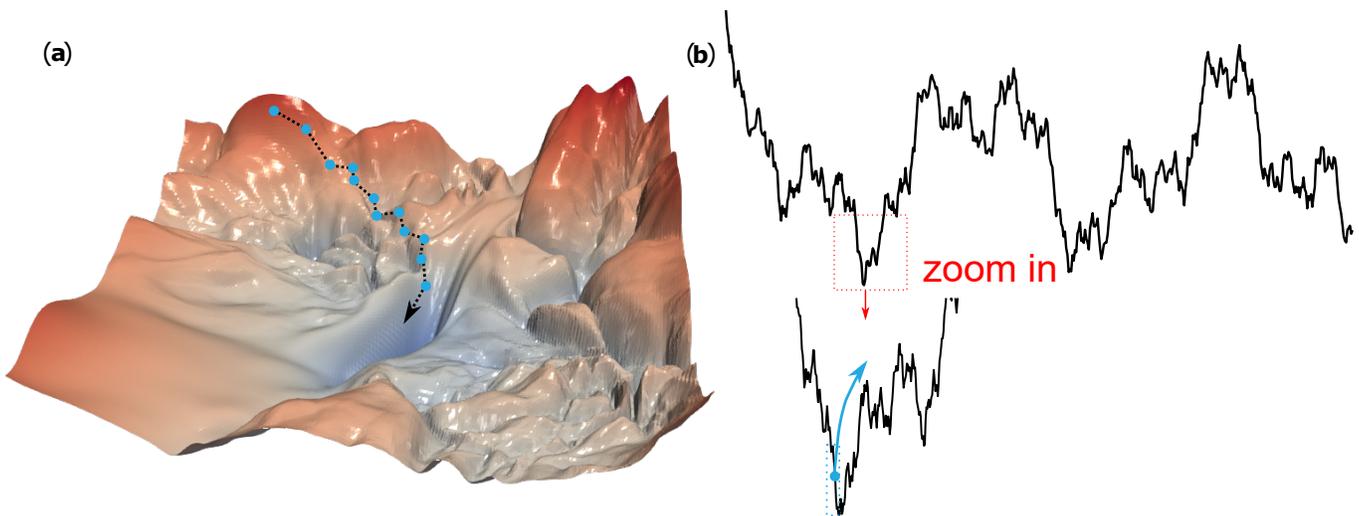


FIG. 1. **Schematic diagram of non-convex loss landscape with a fractal-like structure.** (a) The log-value of loss landscape projected to 2D shows complex structures. The training process based on SGD can be regarded as a SGD optimizer moving on the loss landscape. (b) The fine structure of non-convex loss landscape with a fractal structure shows self-affine and hierarchical properties.

II. DNNs SETUP AND CHARACTERIZATIONS OF ANOMALOUS DIFFUSION DYNAMICS

A. DNNs setup

We consider two classes of neural networks: 1) ResNet-14/20/56/110 [15], where each type is labeled with the total number of layers it has. 2) “VGG-like” networks that do not contain shortcut/skip connections. We produce these networks simply by removing the shortcut connections from ResNets, termed ResNet-14/20/56/110-noshort. All models are trained by vanilla SGD on multiple datasets including MNIST and CIFAR-10, by using two types of loss functions (i.e., cross-entropy, and Kullback Leibler divergence losses). The training processes each run for 500 epochs. All networks are initialized in the standard procedures of the PyTorch library (version 1.3.0). Source code is available at <https://github.com/ifgovh/Anomalous-diffusion-dynamics-of-SGD.git>.

B. Characterizations of anomalous diffusion learning dynamics

Given that full-batch gradient descent is computationally expensive, in the SGD algorithm, the weight parameters $\mathbf{w} = (w_1, w_2, \dots, w_d)$ are estimated by minimizing the minibatch loss function $\nabla \tilde{L}_t(\mathbf{w}) : \mathbb{R}^d \rightarrow \mathbb{R}$, according to

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \eta \nabla \tilde{L}_t(\mathbf{w}_t), \quad (1)$$

where \mathbf{w}_t denotes the d -dimension weight vector (w_1, w_2, \dots, w_d) at time t , and η is the learning rate. The

partial absence of the dataset generates gradient noise $U_t \triangleq \nabla \tilde{L}_t(\mathbf{w}_t) - \nabla L_t(\mathbf{w}_t)$. The updating rule can be rewritten as:

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \eta \nabla L_t(\mathbf{w}_t) + U_t, \quad (2)$$

Hence the SGD training process is also a random diffusive process, where \mathbf{w} can be geometrically interpreted as coordinates of the SGD optimizer in the loss landscape L (Fig. 1).

The loss function of DNN exhibits complex structures, as demonstrated by the projected 2D loss landscape of ResNet-56-noshort [7] (Fig. 1), analogous to complex energy landscape in physical systems [16–19]. In these physical systems, energy landscapes possess fractal-like structures and anomalous diffusion motions of particles stem directly from such kinds of structures. In this study, we apply similar methods used in these systems to quantify the diffusion dynamics of the SGD optimizer. Particularly, the time-averaged mean squared displacement (MSD) [20, 21] is used to characterize the dynamics of SGD moving through the loss landscape, which is defined as

$$\Delta r^2(t_w, \tau) = \frac{1}{T} \sum_{t=t_w}^{t_w+T} \sum_{i=1}^d (w_i(t+\tau) - w_i(t))^2, \quad (3)$$

where τ is lag time, T is the length of the time interval $[t_w, t_w + T]$, and $w_i(t)$ is the value of the i^{th} weight at time t . t_w is the time lapse after the start of the training process (i.e., waiting time). The time variables t , t_w , and T are in units of iteration, and the unit time step corresponds to a single update of the weights.

We characterize the diffusion dynamics based on time-averaged MSD. Although ensemble-averaged MSD is the

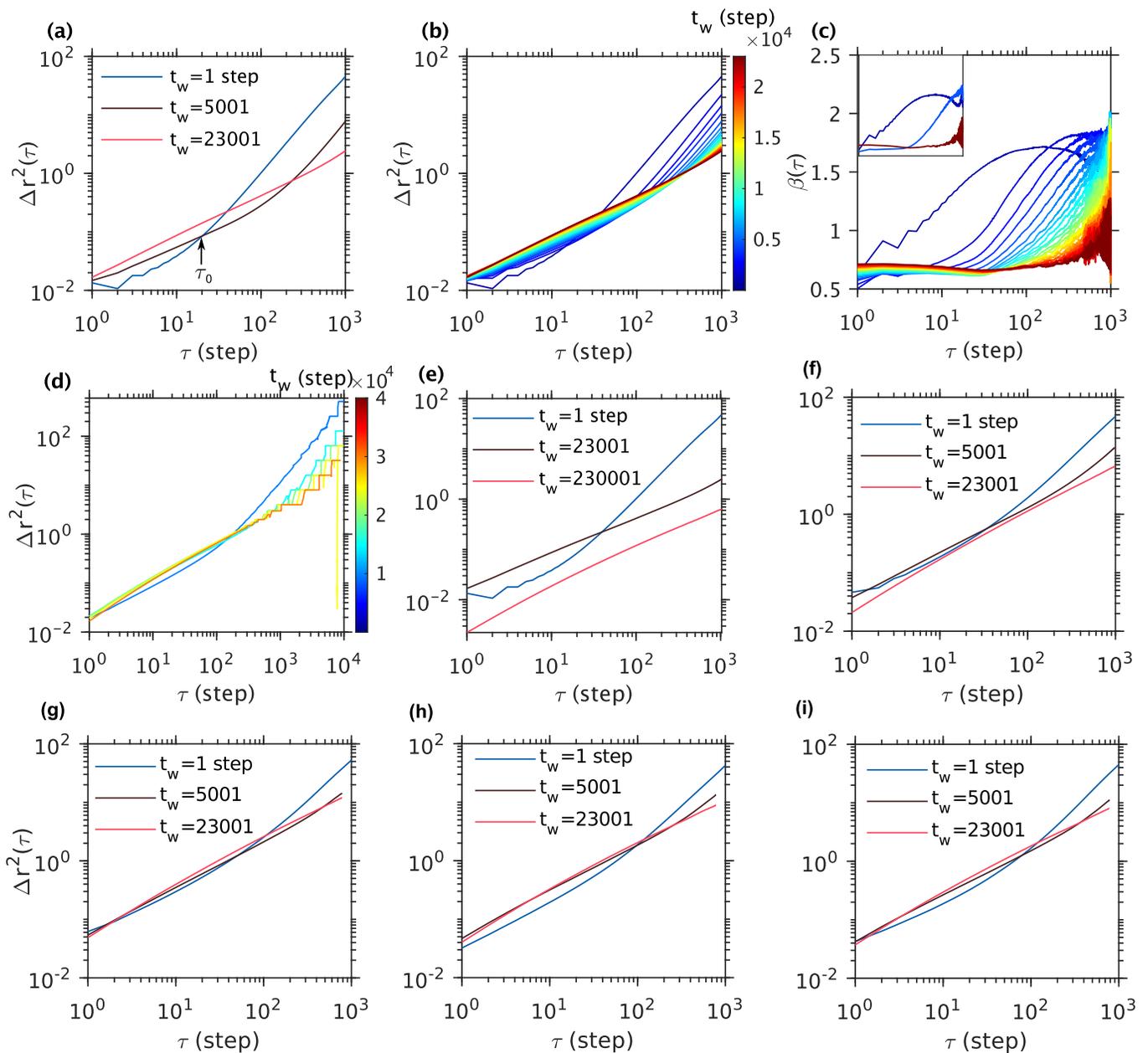


FIG. 2. **SGD dynamics are anomalous diffusion.** (a) MSD $\Delta r^2(\tau)$ of SGD as a function of lag time τ in interval $[t_w, t_w + T]$, $T = 1000$ for ResNet-14 (batch size of 1024, 500 epochs). (b) Same as in (a) but the starting time of the interval t_w gradually increases from 1 to 24000 steps, covering the whole training process. (c) Logarithmic derivative $\beta(\tau)$ of the MSD shown in (b). Inset: Logarithmic derivative $\beta(\tau)$ of the MSD shown in (a). The color scheme is the same as in (b). (d) Same as in (b) but for $T = 10000$. (e) Same as in (a) but for training ResNet-14 5000 epoch. (f-i) Same as in (a) but for ResNet-110 (batch size: 512), ResNet-56 (batch size: 128), ResNet-20-noshort (batch size: 128), and ResNet-20 (batch size: 128), respectively.

most appropriate in theory [22], averaging the ensemble of a high-dimension system is impossible in practice. Therefore, time-averaged MSD is a common tool used to quantify anomalous diffusive processes [20, 21]. No-averaged MSD (Eq. 4) has been used to demonstrate how much the configuration of DNNs and spin glass models

at time $t_w + \tau$ decorrelates from the one at time t_w [11].

$$\Delta(t_w, t_w + \tau) = \frac{1}{d} \sum_{i=1}^d (w_i(t_w + \tau) - w_i(t_w))^2. \quad (4)$$

However, we have calculated no-averaged MSD and found that no-averaged MSD curves are too noisy to characterize the anomalous diffusion learning dynamics.

Brownian motion is identified by $\Delta r^2(\tau) \propto \tau^\alpha$, for large T with the diffusion exponent $\alpha = 1$; the MSD is a linear function of lag time τ . When $\alpha \neq 1$, the corresponding diffusion process has a nonlinear relationship with respect to τ and is defined as anomalous diffusion [23]. If $1 < \alpha < 2$, it is a superdiffusive process; superdiffusion consists of small movements that are intermittently interrupted by long-range jumps. This process has been widely observed in complex physical and biological systems [22]; the mixture of small movements and big jumps in this process is essential for optimally transporting energy in turbulent fluid [24], and for animals to optimally search for spatially distributed food [25]. If $0 < \alpha < 1$, the optimizer is subdiffusive, indicating that it moves slower on average than a normal diffusion process.

III. RESULTS

A. Anomalous diffusion of SGD dynamics

We first illustrate that anomalous diffusion processes characterize SGD learning dynamics. As the findings of different DNN settings are similar, we thus demonstrate all results in ResNet-14 with a batch size of 1024, trained on CIFAR-10 with the cross-entropy loss function and a learning rate of 0.1, unless stated otherwise.

The MSD of each interval $[t_w, t_w + T]$ is calculated according to Eq. 3 ($T = 1000$) for a given t_w . To demonstrate how the diffusion dynamics of SGD optimizer change during the training process, we change t_w systematically. As shown in Fig. 2(a), there are distinct regimes of the learning dynamics characterized by the diffusion exponent α . For the first regime $t_w < t_0$, $t_0 = 21000$ (blue curve, Fig. 2(a)), MSD curves have two segments on the scale $\tau \in [1, T]$ with the smooth transitions around τ_0 (τ_0 is labeled in Fig. 2(a)). When the lag time $\tau > \tau_0$, the MSD curves can be fitted to $\Delta r^2(\tau) \propto \tau^\alpha$ with $\alpha > 1$, indicating that the SGD optimizer exhibits superdiffusive dynamics. However, when $\tau < \tau_0$, $\alpha < 1$, i.e. the SGD dynamics are subdiffusive (Fig. 2(a)). The diffusion exponent α is calculated via the least-squared fitting method. We attempt to fit $\Delta r^2(t_w, \tau) \propto \tau^\alpha$ for $\tau \in [\tau', T]$. We determine $\tau' \in [0, T]$ as the smallest value such that the root-mean-square deviation RMSE < 0.03 ; this τ' is denoted as τ_0 .

During the initial phase of the training process, the interval of the superdiffusion is much longer than that of the subdiffusion. Nevertheless, as t_w increases, the superdiffusion gradually attenuates, as demonstrated by the decrease of the diffusion exponent α and the increase of τ_0 (brown curve in Fig. 2(a); all curves are shown in Fig. 2(b)). In the second regime $t_w > t_0$, the diffusion exponent $\alpha = 0.78$, i.e., the subdiffusion process eventually becomes dominant, as shown by the red curve in Fig. 2(a). To identify the change from the first regime to the second one, we estimate t_0 by fitting the MSD curves

($[t_w, t_w + T]$). To do this, we shift t_w from 1 to 23001, and $[t_0, t_0 + T]$ is the first curve whose goodness of fit has RMSE < 0.03 (0.03 is the standard deviation of all RMSE values). These phenomena can be summarized as below:

$$\Delta r^2(t_w, \tau) \propto \begin{cases} \tau^{\alpha_1} & \text{if } t_w < t_0, \tau < \tau_0 \\ \tau^{\alpha_2} & \text{if } t_w < t_0, \tau \geq \tau_0 \\ \tau^{\alpha_3} & \text{if } t_w \geq t_0 \end{cases} \quad (5)$$

where $\alpha_1, \alpha_3 \in (0, 1)$ and $\alpha_2 > 1$. Such complex dynamics can be further quantitatively characterized by the logarithmic derivative of the MSD, β [26, 27],

$$\beta(t_w, \tau) = \frac{\ln \Delta r^2(t_w, \tau + \Delta\tau) - \ln \Delta r^2(t_w, \tau)}{\ln(\tau + \Delta\tau) - \ln \tau}, \quad (6)$$

where ($\Delta\tau = 20$, Fig. 2(c)); in the first regime, β quickly increases from a value smaller than 1 to a value larger than 1, but in the second regime, β is larger than 1 only when $\tau > 532$.

The time-inhomogeneous anomalous diffusion dynamics are not sensitive to the interval T . For other values such as $T = 5000$, $T = 10000$, the SGD process exhibits similar dynamics with a change from a superdiffusion-dominated regime to a subdiffusion one. Figure 2(d) shows an example of $T = 10000$; initially, the SGD optimizer shows subdiffusion when $\tau < \tau_0$ (blue curve); gradually, the superdiffusion attenuates and subdiffusion process eventually emerges (red curve). Note that we show the results in 500 epochs (24000 steps) because the dynamics after 500 epochs remain the same. The same model is also trained for 5000 epochs, the corresponding MSD curves for $t_w > 24000$ still demonstrate subdiffusion (Fig. 2(e)).

In addition, the time-inhomogeneous anomalous dynamics are not specifically unique to DNN models. Figures 2(f-i) illustrate several models and they also demonstrate similar learning dynamics, including ResNet-14 with the batch size of 128, ResNet-14-noshort with the batch size of 1024, and ResNet-56 with the batch size of 1024. These models are trained with a learning rate of 0.1. This result thus indicates that the SGD learning dynamics generally possess an initial superdiffusion-dominated phase, which gradually evolves to a subdiffusion phase.

The anomalous diffusion learning dynamics provide a way to characterize contributions of network structures and different hyperparameters to the training process. To demonstrate this, we train two types of DNNs, i.e., ResNet and VGG-like networks. VGG-like networks are produced simply by removing shortcut connections from ResNets. As shown in Figs. 3(a-b) and (d-b), shortcut connections do not change the diffusion exponent α significantly in both regimes, with α being greater than 1. However, they do affect the scale range of superdiffusion characterized by τ_0 when $t_w = 1$. τ_0 indicates the crossover of MSD in the first 1000 steps. Specifically,

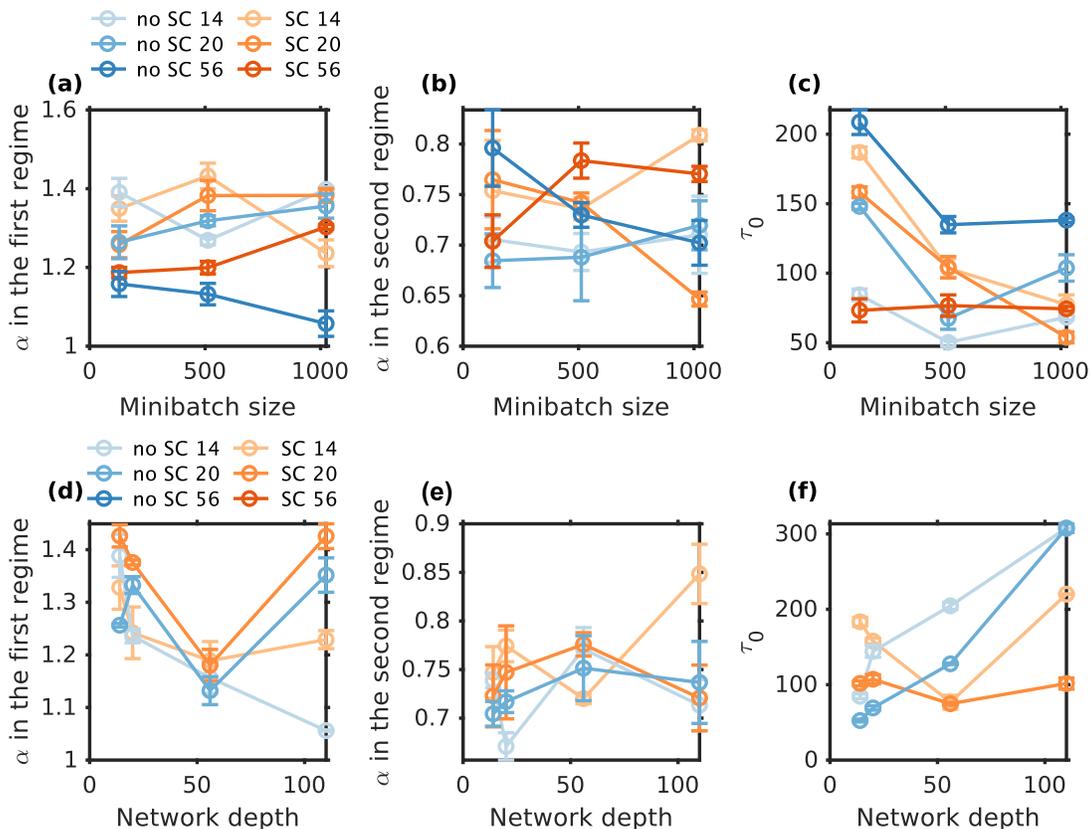


FIG. 3. **The effects of the depth, batch size, and shortcut connections of DNNs on the anomalous diffusion dynamics of SGD.** The orange and blue colors represent the network with/without shortcut connections (SC) respectively. The digits in legends of the first row (a-c) represent network depth; for example, no SC 14 denotes ResNet-14-noshort. Those in the second row (d-f) represent minibatch size; for example, no SC 128 denotes ResNet training using the minibatch size of 128. (a) The diffusion exponents α on larger lag times ($\tau > \tau_0$) when $t_w = 1$ as a function of minibatch size. (b) Similar to (a) but for α in the second regime (when $t_w = t_0$). (c) The crossover τ_0 as a function of minibatch size. Here, τ_0 is the lag time when the MSD curve transitions from subdiffusion to superdiffusion when $t_w = 1$. (d-f) Similar to (a-c) but for varying network depth. The error bars represent the standard deviation calculated over 10 trials.

with shortcut connections, τ_0 is smaller than those without shortcut connections, indicating that the scale of superdiffusion is elongated (Fig. 3(c) and Fig. 3(f)). The superdiffusion dynamics enable the SGD walker to explore larger areas of loss landscape in a fixed time than normal diffusion and subdiffusion processes. From this perspective, DNNs with shortcut connections can facilitate training, which is consistent with previous studies [7, 15].

The anomalous diffusion learning dynamics are not very sensitive to minibatch sizes. As shown in Fig. 3(a), the diffusion exponent α in the first regime does not significantly vary ($\pm 15\%$) with respect to the change of minibatch size from 128 to 1024 in all networks, although τ_0 decreases in ResNet-14,20 with the increase of minibatch size (Fig. 3(c)).

Furthermore, as shown in Fig. 3(d), α changes nonlinearly with respect to the network depth. However, the network depth reduces the scale range of superdiffusion, characterized by τ_0 (Fig. 3(f)); this result suggests that the deeper DNNs are more difficult to be trained [28] due

to the shorter scale of superdiffusion.

We next study the effect of the learning rate η on the anomalously diffusive learning dynamics. By varying learning rates from 0.001 to 0.5, we find that it only influences the emerging sequence of the superdiffusion learning dynamics (Fig. 4(a)). As shown in Fig. 4(b), DNN training with small learning rates is much slower than that with large learning rates. Thus, a small or large learning rate can only slow down or speed up the training procedure, but does not change the fundamental occurrence of anomalous diffusion dynamics. When the learning rate is small ($\eta < 0.05$), during 500 epochs, the training processes remain in the first regime; there is no pure subdiffusion in Fig. 4(c) and Fig. 4(d). For example, when $\eta = 0.001$, the MSD curves for $t_w < 5000$ have only one segment instead of two and a diffusion exponent of $\alpha = 2$. This is because small learning rates delay the training process. However, DNNs trained with larger learning rates grant superdiffusion in the first regime and pure subdiffusion dynamics in the second regime, as shown in Fig. 4(e) and Fig. 4(f).

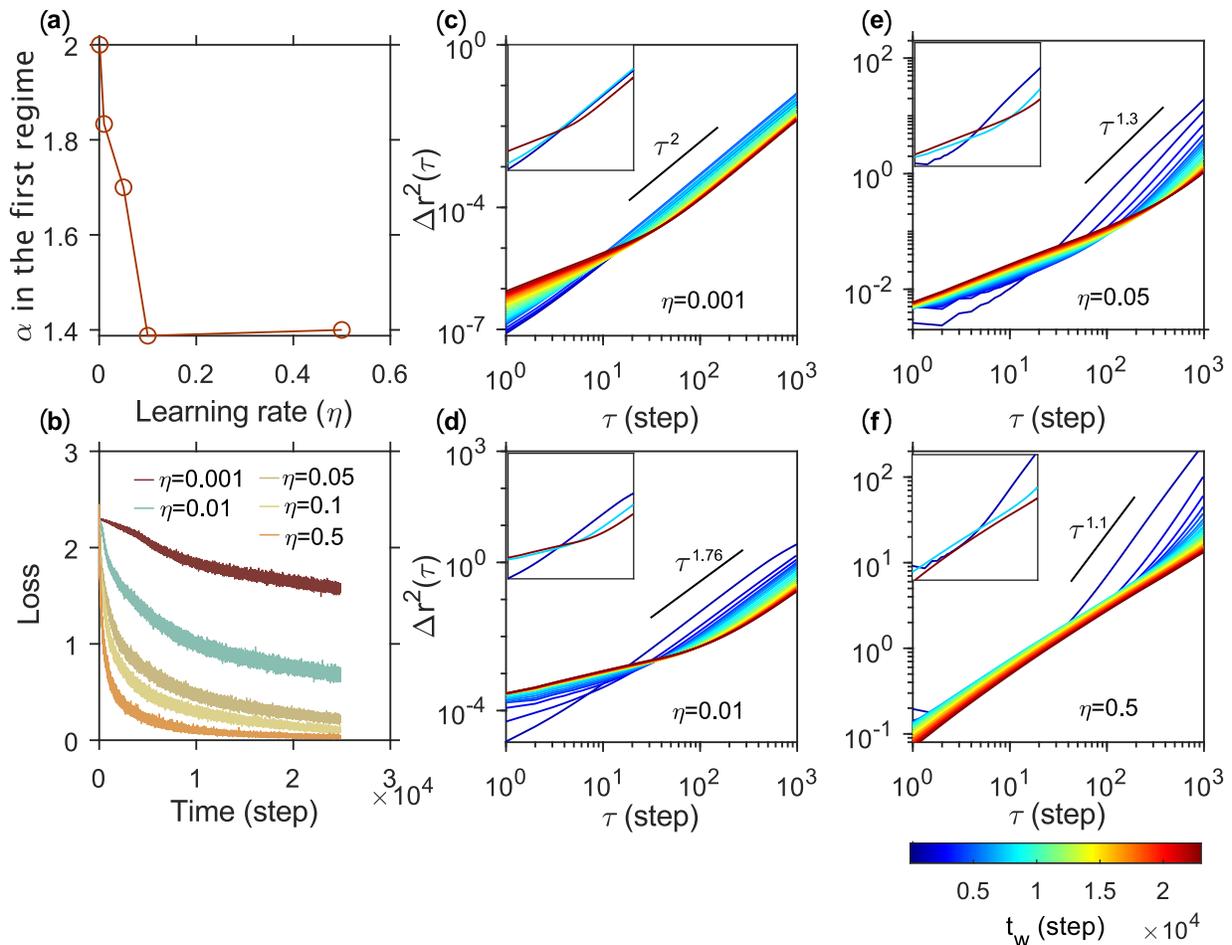


FIG. 4. **The effect of learning rate on the anomalous diffusion dynamics of SGD.** All results are from ResNet-14 with a batch size of 1024 on CIFAR10 dataset. (a) The diffusion exponent α on larger lag times ($\tau > \tau_0$) when $t_w = 1$ as a function of learning rate η . (b) The loss value as a function of time. (c-f) The MSD of SGD for the learning rates of 0.001, 0.01, 0.05, and 0.5, respectively. Jet colormap represents the starting time point (t_w) as in Fig. 2(b). One curve represents MSD in 1000 steps. Inset: The MSD of SGD when $t_w = 1, 7001, 23001$. The black lines are eye guides.

B. Heavy-tailed gradients

To gain further insights into the physical origin of the anomalous diffusion dynamics, we next demonstrate the statistical property of minibatch gradients $\nabla\tilde{L}$. We first introduce the definition of the Lévy α -stable distribution. Given a Lévy stable random variable X , it is characterized by the characteristic function [29]

$$\varphi(u; \alpha_{\text{dist}}, \beta, \gamma, \delta) = \exp(iu\delta - |\gamma u|^{\alpha_{\text{dist}}}(1 - i\beta \text{sgn}(u)\Phi)) \quad (7)$$

where $\text{sgn}(u)$ is the sign of u and

$$\Phi = \begin{cases} \tan\left(\frac{\pi\alpha_{\text{dist}}}{2}\right) & \alpha_{\text{dist}} \neq 1 \\ -\frac{2}{\pi} \log|t| & \alpha_{\text{dist}} = 1 \end{cases}$$

α_{dist} is the stability parameter with the range $0 < \alpha_{\text{dist}} < 2$. The probability density function (PDF) decays with a power-law tail $|x|^{-\alpha_{\text{dist}}-1}$ which is slower compared to Gaussian distributions; thus the distribution is heavy-

tailed [30]. When $\alpha_{\text{dist}} = 2$, the distribution is Gaussian. β is the skewness parameter. In particular, for a symmetric Lévy α -stable ($\mathcal{S}\alpha\mathcal{S}$) random variable X , i.e. $X \sim \mathcal{S}\alpha\mathcal{S}$, the skewness parameter $\beta = 0$ which indicates the PDF is symmetric around 0. γ is the scale parameter and δ is the shift parameter.

We next estimate the minibatch gradient $\nabla\tilde{L}$ in Eq. 1 (vanilla SGD) with respect to each w_i at each time point. In the first regime ($t_w < t_0, t_0 = 21000$), it can be fitted to a symmetric Lévy α -stable distribution by the maximum likelihood method; stability parameter $\alpha_{\text{dist}} = 1.46528$ [1.46509, 1.46546] (the brackets denote the 95% confidence interval).

The power-law tail of the distribution of $\nabla\tilde{L}$ shown in Fig. 5(a) inset further validates the heavy-tailed distribution. We further compare log-likelihood ratios between the fitted Lévy α -stable distribution and Gaussian distribution, and find that the log-likelihood ratios (1.52×10^9) are sufficiently positive, indicating that the distributions most likely follow Lévy α -stable distribution ($p < 10^{-15}$,

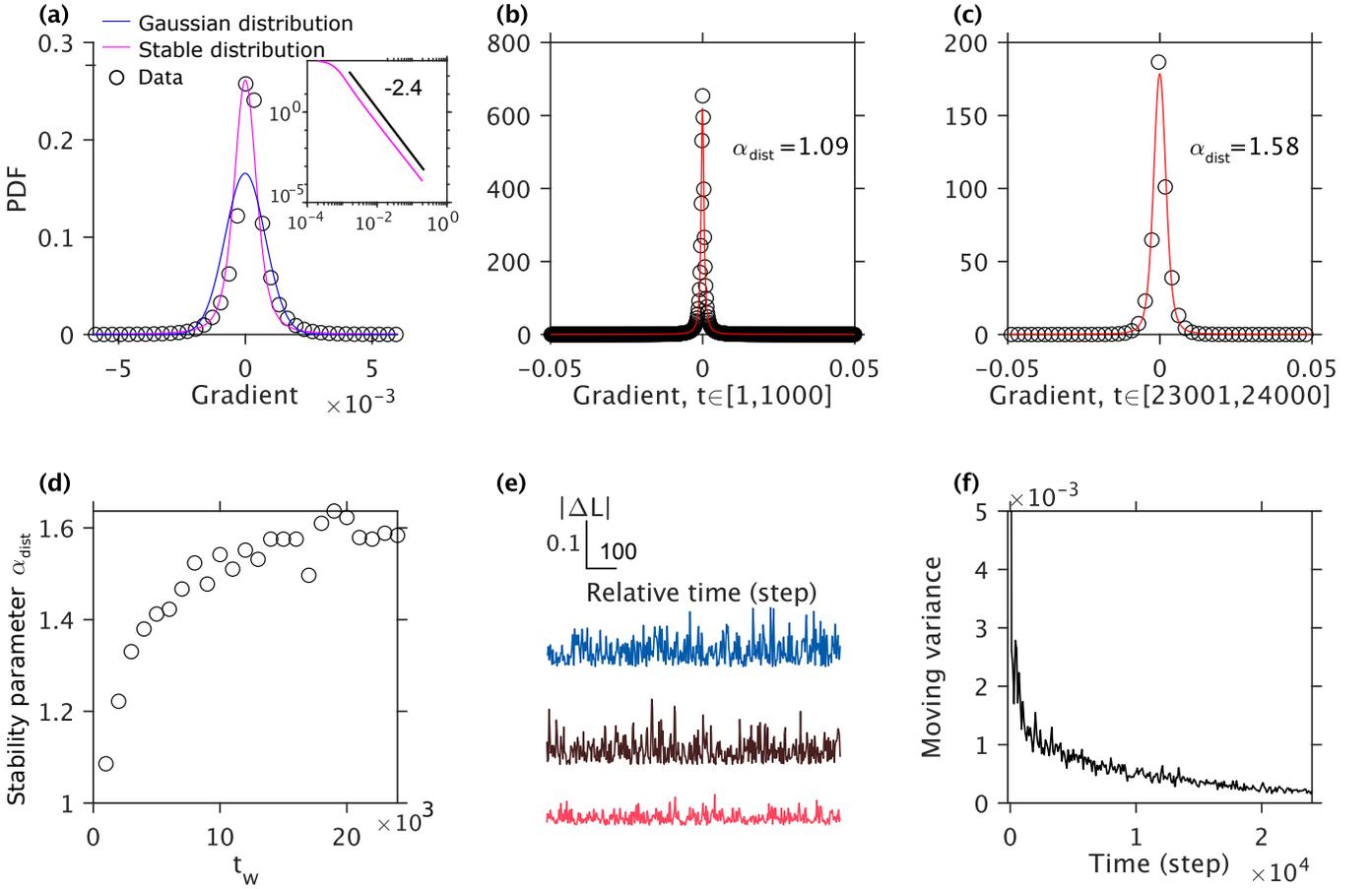


FIG. 5. **The gradient is heavy-tailed.** The distribution of minibatch gradients $\nabla\tilde{L}(\mathbf{w})$ in the first regime can be fitted as a symmetric Lévy α -stable distribution (red curve). Inset: the log-log plot of the positive part of the distribution indicates that it has a power-law tail. (b) When $t \in [1, 1000]$, $\nabla\tilde{L}(\mathbf{w})$ can be fitted in symmetric Lévy α -stable distribution with the stability parameter $\alpha_{\text{dist}} = 1.09$. (c) Same as in (b) but for $t \in [23001, 24000]$. (d) The fitted stability parameter in the interval of $[t_w, t_w + T]$ as a function of t_w . (e) Fluctuations of absolute value of change-of-loss ΔL decrease as t_w increases. Colormap is the same as in Fig. 2(a). (f) The moving variance of loss as a function of time.

Vuong test).

To illustrate the evolution of gradient distribution, we also fit the distributions of gradients in intervals $[t_w, t_w + T]$ to Lévy α -stable distribution; the log-likelihood ratios compared with Gaussian distribution are sufficiently positive. Figure 5(b) demonstrates the distribution in the first interval, $t_w = 1$; the stability parameter $\alpha_{\text{dist}} = 1.09$ [1.05, 1.12]. Figure 5(c) demonstrates the distribution in the last interval, $t_w = 23001$ and correspondingly $\alpha_{\text{dist}} = 1.58$ [1.54, 1.63]. The stability parameter α_{dist} increases as training evolves (Fig. 5(d)), indicating a reduction in the heavy-tailedness of gradient distribution. Because the heavier the tail, the larger the fluctuations of gradient values, and as the changes of gradient are directly related to the MSD (Eq. 1 and Eq. 2), the result regarding the changes of gradient distributions is consistent with the time-inhomogeneous anomalous diffusion dynamics where superdiffusion attenuates gradually to subdiffusion. It is also interesting to note that in the physics literature, it has been found that the increase

of the heavy-tailedness of step sizes of random walkers results in super-diffusive motions [22]. Such superdiffusive processes with intermittent long-range jumps might help the optimizer jump out local minima, facilitating basin hopping during the initial exploratory phase of training. This point is further illustrated below, based on a simple model of SGD. Entering and leaving local minima give rise to the fluctuations of loss values (L). To demonstrate the changes of these fluctuations, we first calculate the absolute value of change-of-loss $|\Delta L| = L(\mathbf{w}(t+1) - \mathbf{w}(t))$. As shown in Fig. 5(e), as t_w increases, $|\Delta L|$ decreases. Such behavior is further quantified by the decreasing moving variance of the loss L against time (Fig. 5(f)); the moving variance is calculated over a sliding window of 100 steps across neighboring L .

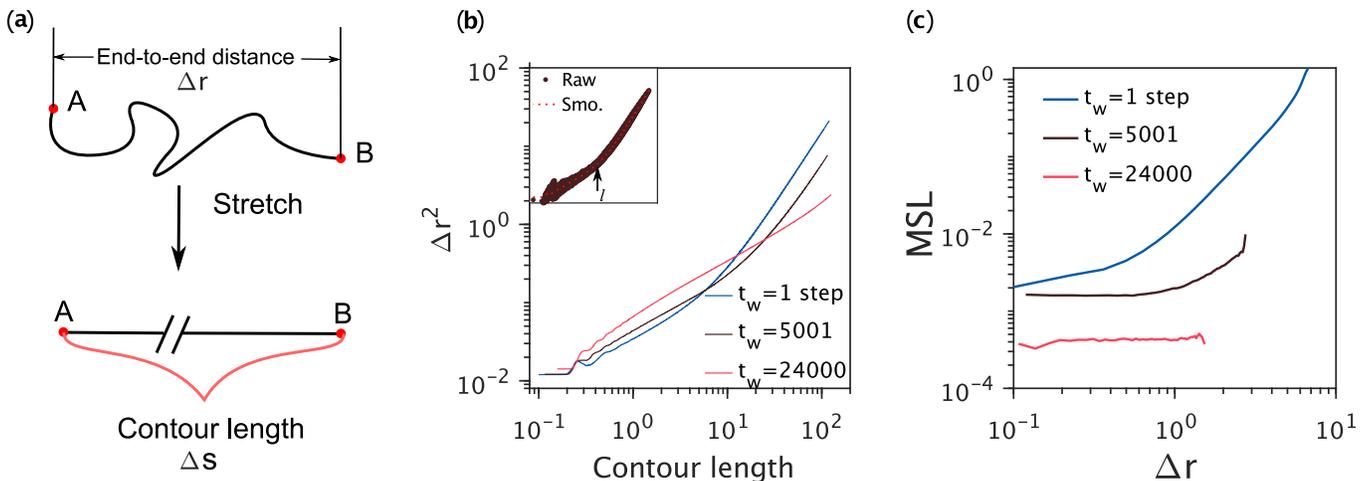


FIG. 6. **The fractal-like loss landscape of DNN and SGD path.** (a) Schematic diagram of contour length Δs and end-to-end distance Δr . (b) Squared end-to-end displacement (Δr^2) as a function of contour length. The curve is smoothed by moving average over each window of raw data for clearer illustration; the window length is 40 steps. The inset displays the raw data when $t_w = 1$ and the dashed red line represents the smoothed data. (c) Mean squared loss differences of point pairs (MSL) against the end-to-end distances Δr in the interval $[t_w, t_w + T]$. The colormaps in (b) and (c) are the same.

C. Fractal trajectory of SGD

In complex physical systems, the anomalous diffusion of particles may associate with its fractal trajectories [17]. To verify that the SGD path is fractal, we characterize the power-law scaling of its end-to-end length and contour length as used in [17]. Specifically, the contour length is the path length from one end to another and is calculated by accumulating step sizes (curve length) as illustrated in Fig. 6(a). As shown in Fig. 6(b), the segmented end-to-end distance (Δr) of SGD path appears to scale with its contour length, Δs . The data is smoothed to be clear; the example of raw data and the smoothed data are shown in Fig. 6(b) inset. The blue and brown curves in Fig. 6(b) display two distinct scaling regimes, illustrating the fractal property of SGD paths in certain ranges [17]. The fractal dimension D_f is calculated by the scaling of Δs and Δr ($\Delta r^2 \sim \Delta s^\lambda$), $D_f = (2/\lambda) \in [1.32, 2.67]$. Separated by the crossover l (labeled in Fig. 6(b) inset), λ on longer length scales ($\Delta s > l$) is larger than that on short length scales ($\Delta s < l$). As t_w increases, MSL collapse to a power law with an exponent of 0.75 (red curve in Fig. 6(b)). The time-inhomogeneous dynamical changes of SGD trajectory are similar to the case of the MSD (Fig. 2(a)).

To further determine whether the SGD path is self-affine or self-similar, we calculate the transverse distance and compare it with the end-to-end distance. The transverse distance between two points on the trajectory is the maximal distance perpendicular to the straight line connecting these two points [17]. We find the transverse distances of different points along paths do not scale to their end-to-end distances; this indicates self-affine rather than self-similar fractal, because the former contains non-uniform scaling, i.e. the shapes are (statisti-

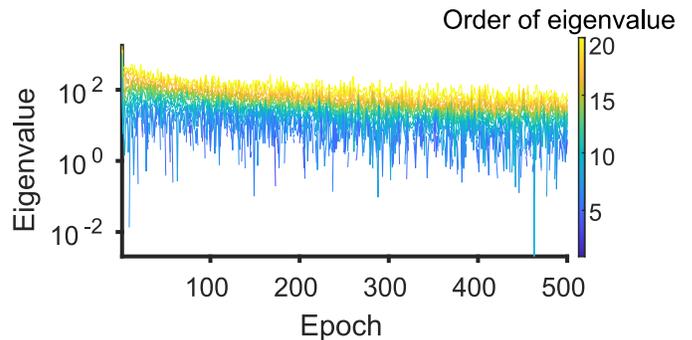


FIG. 7. **The top 20 eigenvalues of Hessian matrix of loss landscape decrease as training epoch increases.** This indicates that the SGD optimizer leaves narrow local minima and enters a flatter area in the loss landscape of DNN.

cally) invariant under transformations that scale different coordinates by different amounts [31].

D. Fractal-like loss landscape

In physical systems, anomalous diffusion motions and fractal path of moving particles could be consequences of the energy landscape itself being fractal [16, 17, 19]. Inspired by these studies, we hypothesize that the loss landscapes of DNNs could have fractal-like structures. It is impossible to directly quantify the fractal dimensions of high-dimensional loss landscape; we thus use an approach proposed in [32]. From [31], the definition of fractal is as the following:

Definition A random function f on a metric space is fractal if the distribution of $f(X')$ conditional on $f(X)$,

X , and X' , is normal with mean 0 and variance proportional to $\Delta r(X', X)^{2H}$, where H is a parameter in $(0, 1)$ and $\Delta r(X', X)$ is the distance between point X and X' .

The distribution of X and X' in the definition is over the probability space from which f is drawn, but experimentally we sample over random values of X and X' for a given sample function f . The equivalence of these two procedures is referred to as ergodicity and we take it for granted. Checking the distributions for each value of distance in practice is difficult; hence we measure the expectation of $[f(X') - f(X)]^2$ to characterize the fractal-like structure and check if it satisfies Eq. 8. In the context of a SGD optimizer on the loss landscape, f is the loss function L ; X and X' are a pair of weights (\mathbf{w} and $\tilde{\mathbf{w}}$) on the loss landscape. Thus, the expectation is referred to as mean squared loss (MSL),

$$\text{MSL} \propto \Delta r(\mathbf{w}, \tilde{\mathbf{w}})^{2H} \quad (8)$$

where $\Delta r(\mathbf{w}, \tilde{\mathbf{w}})^{2H}$ is the end-to-end distance between \mathbf{w} and $\tilde{\mathbf{w}}$. Note that in the field of machine learning, the same method has been used to quantify fractal landscapes for simulated annealing [32].

The MSL is calculated by the following equation,

$$\text{MSL}(t_w, \Delta r) = \frac{1}{N_{\Delta r}} \sum_{j=1}^{N_{\Delta r}} [L(\mathbf{w}_{\Delta r}^{t_w, T}(j)) - L(\tilde{\mathbf{w}}_{\Delta r}^{t_w, T}(j))]^2, \quad (9)$$

where the pair of weights $\mathbf{w}_{\Delta r}^{t_w, T}(j)$ and $\tilde{\mathbf{w}}_{\Delta r}^{t_w, T}(j)$ are sampled at different time steps along the trajectory in $[t_w, t_w + T]$ with the end-to-end distance between them equal to Δr , and $N_{\Delta r}$ is the total number of pairs at a distance of Δr . To be consistent with piecewise MSD, we choose $T = 1000$ steps. We use the points sampled by SGD (i.e., points along the optimization trajectory) to estimate MSL; it illuminates certain structure on the local area of the loss landscape that we would like to explore.

As shown in Fig. 6(c) (blue curve), in the first regime ($t_w < t_0$), the MSL curve can be fitted to a power-law function with an exponent of 1.8 on the larger distance scale ($\Delta r \in [0.4, 10]$). It satisfies Eq. 8 and indicates that the loss landscape of DNN has fractal-like structures at the initial phase of learning process. Note that the power-law scalings do not hold within the whole scale ($[0.1, 10]$). As the superdiffusion attenuates, the end-to-end distance Δr in T decreases and the power-law exponent of the MSL with respect to Δr flattens from period to period (brown curve in Fig. 6(c)). Eventually, in the second regime ($t_w > t_0$), the MSL is around a constant value against varying Δr (red curve in Fig. 6(c)). Based on the above definition, this indicates that the optimizer now reaches a relatively flat region on the landscape. We use ‘‘flat’’ colloquially to indicate approximate flatness [14, 33–35]. Our results thus indicate that the SGD optimizer moves from rougher (more fractal-like) to relatively flatter regions of the loss landscape. The change to the flatter regions can also be quantitatively demonstrated

by the fact that eigenvalues of the Hessian matrix gradually decrease to near-zero values (Fig. 7), which has also been found in [36]. Note that there has been increasing evidence showing that the optimizer can eventually find a good generalizable solution existing at the flat regions of the loss landscape [33, 35, 37]. Our work suggests that it would be relevant to use the methods in these previous studies to characterize how the SGD moves from rougher to flatter regions of the loss landscape.

Importantly, we find that when $t_w < t_0$, the fractal-like hierarchical structure provides highly fluctuating gradients and thus superdiffusion emerges; when $t_w > t_0$, the flatter structure causes fewer fluctuations of gradients and results in subdiffusion. In a recent study [14], it has been shown that during the training process of DNNs, the SGD moves towards flatter regions of the loss landscape and correspondingly the anisotropic SGD noise strength decreases. In future work, it would be interesting to explore whether the changes from superdiffusion to subdiffusion underlie the changes of noise strength as reported in [14].

To further justify the fractal-like structure of loss landscape, we use the method of filter-wise normalized directions [7] to project the loss landscape of ResNet-14 (batch size of 1024, learning rate of 0.1, trained on CIFAR-10 with cross-entropy loss function) to 2D space. We then calculate the Hausdorff dimension of a 2D projected loss landscape via the box-counting method [38]; the fractal dimension is approximately 1.8.

E. Fractal-like landscapes can cause anomalous diffusion learning dynamics

To further understand the contributions of fractal-like loss landscapes to the anomalous diffusion of learning dynamics, we develop a simplified model of SGD with a 2D fractal loss landscape, as described below.

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \eta \nabla L(\mathbf{w}_t) + \eta \sigma Z_t, \quad (10)$$

where \mathbf{w}_t is the weight parameters at time t , η is the learning rate, and Z_t is drawn from Gaussian distribution $\mathcal{N}(\mathbf{0}, \sigma)$. The landscapes of L are fractal on a 2D space and are generated by the algorithm in [39] with fractal dimension $\in [1, 2]$ (Fig. 8(a)). Among the generated landscapes, those which contain a wide minimum are selected, simulating the flat basin (local optima) in loss landscapes of DNNs. In such landscapes, the SGD walker moves from rougher regions to flatter regions of the loss landscape, as observed in DNNs. We iterate Eq. 10 1000 times with the learning rate $\eta \in [0.002, 0.01]$ and $\sigma \in [0.005, 0.05]$, and the gradient ∇L is calculated by the numerical gradient of the landscape. The SGD optimizer moves from a high-altitude point to the global minimum, as the example ($\eta = 0.02$ and $\sigma = 0.01$) shown in Fig. 8(a). The MSD also illustrates that the SGD optimizer is dominantly superdiffusive when $t_w < t_0 = 139$ (before entering the final minimum; curve with a square

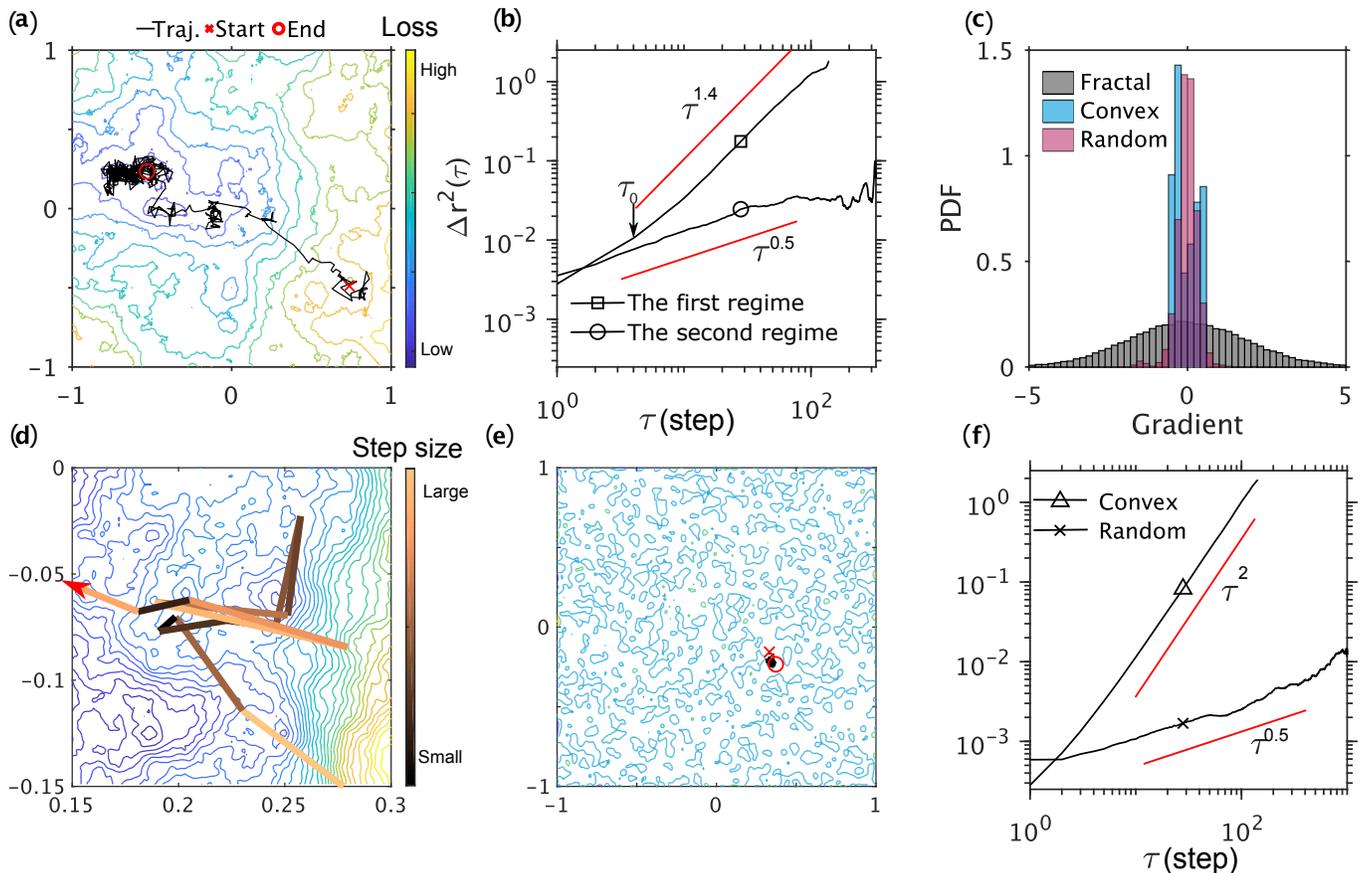


FIG. 8. **The simplified model unravels the anomalous diffusion nature of the SGD optimizer.** (a) Black curve represents the trajectory of the SGD optimizer in the simplified model on the 2D fractal loss landscape illustrated by the contour plot. (b) MSD of the SGD optimizers in the fractal landscape as the function of lag time. The curve is averaged over 10 trials. (c) Distributions of gradients (ΔL) in fractal, convex, and randomly shuffled landscapes, respectively. (d) Trajectory of the SGD optimizer and the landscape in (a) are zoomed in around a local minimum. The colormap from black to golden encodes the step sizes. (e) Same as in (a) but for random landscape. (f) Same as in (b) but for random landscape and convex landscape.

in Fig. 8(b)). When $t_w > t_0$ (after entering the final minimum), the dynamics of SGD optimizer are only subdiffusive as shown in the curve with a circle in Fig. 8(b). Such behaviors demonstrate that the simple model can reproduce the time-inhomogeneous MSD of SGD as found in DNNs.

In this model with 2D fractal landscape, the fractal landscape generates a heavy-tailed distribution of gradients (∇L) of all steps (Fig. 8(c)). As found in the DNNs, the gradient distribution can be fitted as a symmetric Lévy α -stable distribution which has the stability parameter $\alpha_{\text{dist}} = 1.9$ [1.88749, 1.91298]. The goodness of fit is verified by the positive log-likelihood ratio (39.37) of Lévy α -stable distribution and normal distribution. In our model, ∇L is calculated by the numerical gradient of the fractal landscape when $t_w < t_0$. Such heavy-tailed gradients provide a relatively higher possibility of long-range jumps, which generate superdiffusion. In the fine structure of fractal landscape, there are large gradient values (Fig. 1(b)) which propel the SGD optimizer to

jump out narrow minima. It is important to note that as noise in our simplified SGD model is Gaussian, such heavy-tailed gradients and superdiffusion dynamics solely result from the fractal-like loss landscape.

To explicitly illustrate the benefits of fractal landscapes for facilitating the SGD to jump out local minima, we focus on some regions around local minima. As shown in Fig. 8(d), the optimizer moves to the local minimum where the SGD optimizer displays short-range movements as illustrated by darker bars in Fig. 8(d) and long-range movements as illustrated by lighter bars in Fig. 8(d). Some long-range steps make the SGD optimizer escape the minimum; however, some of them jump to a lower altitude and then leave the minimum. This example illustrates how the fractal landscape assists the SGD optimizer escape local minima. As we only use Gaussian noise in this simple model, the main source providing long-range jumps to escape local minima is the heavy-tailed gradients (Fig. 8(c)). This result suggests that in DNNs, the similar mechanism might en-

able the SGD walker to escape local minima in the initial superdiffusion-dominated regime.

Although the results in the simple model are largely consistent with DNNs, there are some differences. In the simple model, if the learning rate is small (< 0.002), the toy model becomes sensitive to initial conditions (the initial position of random optimizers). The random optimizer would be trapped in a local minimum. For large learning rates (> 0.01), combining with the occasional long-range jump, the optimizer would easily go out the landscape.

However, if we choose other types of landscapes and maintain η and σ , the complex MSD dynamics no longer hold. When the landscape is generated by smoothing a randomly shuffled fractal landscape with a Gaussian kernel (standard deviation: 8) such that it is at least once-differentiable, the optimizer gets stuck in a local minimum (Fig. 8(e)) and exhibits subdiffusion (Fig. 8(f)). On the other hand, when the landscape is convex, for example, a convex paraboloid, the MSD has an exponent close to 2 (Fig. 8(f)), inconsistent with the results in DNNs. In comparison to the tail of gradient distribution in fractal landscapes, the ranges of gradients in convex or random landscapes are far smaller (Fig. 8(c)). These results thus indicate the fractal-like loss landscape is essential for the emergence of the anomalous diffusion learning dynamics.

IV. DISCUSSION

We have revealed the anomalous diffusion nature of deep learning dynamics which arises from the interactions of the SGD walker with the geometry structure of the loss landscape. Particularly, we have demonstrated that the fractal-like loss landscape can give rise to the superdiffusion learning dynamics with intermittent big jumps during the initial training phase, which plays an essential role in preventing the SGD optimizer from being trapped in narrow minima. Subdiffusion on the other hand also occurs naturally during the final stage of the training process, stabilizing the movement of the optimizer gradually, potentially when wide minima of landscape are encountered. In addition, we have developed a new SGD model to reveal the mechanistic relations between the fractal landscape, the superdiffusive learning dynamics and their computational benefits. Our results reveal the effectiveness of deep learning from the perspective of its rich, complex dynamics and have implications for designing efficient deep neural networks.

Previous studies tackled the problem from modeling the SGD as a process with heavy-tailed behaviors [13, 40, 41]. Particularly, Simsekli et al reported a heavy-tailed behavior in the stochastic gradient noise (U_t in Eq. 2) and proposed modeling the SGD dynamics as a stochastic differential equation driven by an alpha-stable process. They further invoked existing metastability theory to justify why these dynamics would prefer wide minima [13]. The SGD updating rule can be

represented in terms of gradient noise, i.e., $\mathbf{w}_{k+1} = \mathbf{w}_k - \eta \nabla L(\mathbf{w}_k) - \eta \cdot \text{gradient noise}$. However, the distribution of gradient noise was further found to be Gaussian in the early phases, and changes throughout the training process [42]. To model the gradient noise as a single type of noise is thus inappropriate. Rather than the gradient noise, our work focuses on the change of the drift term throughout the training process (i.e., gradient, $\nabla L(\mathbf{w}_k)$) which is directly related to the structure of loss landscape.

The structure of loss landscape can give rise to the anomalous diffusion dynamics of learning process quantified from MSD [17]. In [11], no-averaged MSD was used to analyze learning dynamics and the slope of MSD is not equal to 1. However, this study did not introduce anomalous diffusion in the discussion of learning dynamics but only showed the aging phenomena. Ideally, we would like to obtain the ensemble average of the MSD, but this is unrealistic for DNNs. We instead used the time-averaged MSD [20–22, 27, 43, 44], which can better demonstrate the anomalous diffusion in the context of DNNs.

By using a similar methodology as in simulated annealing [32], we have quantitatively demonstrated that the loss-landscape of DNN is fractal-like. The fractal-like structure can give rise to heavy-tailed gradients which may help the SGD optimizer to jump out local minima. Also, the fractal landscape may result in fractal trajectories of the SGD optimizer, as in complex physical systems [17]. Our results show that the SGD trajectories are indeed fractal as quantified by the contour length and end-to-end length. Recently, it has been suggested that the fractal trajectory of SGD optimizer may facilitate generalization in DNNs [45]. Such fractal trajectories might result from the fractal-like structure of loss landscape as what we have demonstrated. Our results based on both the MSD and MSL indicate that the SGD walker moves from rougher (more fractal-like) areas to flatter areas of the loss landscape. During this process, the learning dynamics change from superdiffusion-dominated dynamics, which help the SGD to escape from local traps, to subdiffusive dynamics which can rather consolidate the residence of the SGD in the flatter areas with good solutions (minima). These time-inhomogeneous anomalous diffusion learning dynamics arising from the interactions of SGD and the loss landscape thus provide insights into understanding how the optimizer can find flat minima.

The simple SGD model on a 2D fractal landscape generates the same pattern of time-inhomogeneous learning dynamics as in the high-dimensional DNNs, which however cannot be accounted for by the traditional formulation based on the Langevin equation with Gaussian noise [8, 9]. The simple SGD model does not involve any type of non-Gaussian noise and demonstrates that fractal landscapes alone can lead to anomalous diffusion learning dynamics, thus indicating that the interactions between the fractal loss landscape and SGD are the mechanism underlying the emergence of the anomalous diffusive learning dynamics. However, as our 2D model is not

directly derived from DNN models, the generalization of this mechanism to DNNs is limited. Anomalous superdiffusion and subdiffusion are nonlinear diffusive processes and are generally referred to as fractional motions that can be formulated based on fractional differential equations [23], suggesting that developing a fractional mean field theory as in [46] for understanding deep neural networks would be a promising direction to pursue in the future. In addition, future studies should figure out the major source of fractal-like loss landscape. The training landscape is composed of the data and the network architecture. Some previous studies have shown that realistic datasets such as handwritten digits (MNIST), rather than random noise, have low-dimension structure/manifold [47–49]. It would be interesting to study the effect of such data structure on the geometrical properties of loss landscape.

On the other hand, the network structure can affect anomalous diffusion learning dynamics. We have found that the deeper DNNs, the shorter scale of superdiffusion, indicating a more demanding training process, consistent with the empirical rules of DNNs [28, 34, 50] and extended the understanding from the aspects of training dynamics and landscape structures [12]. Additionally, shortcut connections in ResNet can extend the scale of superdiffusion, explaining why employing such tech-

niques reduces the difficulties of training DNNs. These findings agree with the theoretical and experimental results of gradient confusion [28] and the visualization of 2D projected loss landscapes [7]. These studies found that shortcut connections reduce the difficulty in the training process by smoothing out the loss landscape. Furthermore, we find that the batch size cannot significantly alter the anomalous diffusion dynamics, supporting the conclusion that gradient noise is not the only driving force to escape critical points (saddle points or local minima). For future studies, it remains important to find out the quantitative relation between other architectures such as batch normalization [50] and fractal-like landscape structure.

ACKNOWLEDGMENTS

The authors acknowledge the University of Sydney HPC service for providing high-performance computing that has contributed to the research results reported within this paper. This work was supported by the Australian Research Council (grant nos. DP160104316, DP160104368).

-
- [1] Y. LeCun, Y. Bengio, and G. Hinton, *Nature* **521**, 436 (2015).
 - [2] Z. C. Lipton, in *The International Conference on Learning Representations (ICLR) workshop* (2016).
 - [3] E. Hoffer, I. Hubara, and D. Soudry, in *Advances in Neural Information Processing Systems* (2017) pp. 1731–1741.
 - [4] T. J. Sejnowski, *Proceedings of the National Academy of Sciences*, 201907373 (2020).
 - [5] A. Choromanska, M. Henaff, M. Mathieu, G. B. Arous, and Y. LeCun, in *Artificial Intelligence and Statistics* (2015) pp. 192–204.
 - [6] S. Becker, Y. Zhang, and A. A. Lee, *Physical Review Letters* **124** (2020).
 - [7] H. Li, Z. Xu, G. Taylor, C. Studer, and T. Goldstein, in *Advances in Neural Information Processing Systems* (2018) pp. 6389–6399.
 - [8] S. Jastrzębski, Z. Kenton, D. Arpit, N. Ballas, A. Fischer, Y. Bengio, and A. Storkey, arXiv:1711.04623 (2017).
 - [9] M. Welling and Y. W. Teh, in *Proceedings of the 28th International Conference on Machine Learning* (2011) pp. 681–688.
 - [10] P. Chaudhari and S. Soatto, in *Information Theory and Applications (ITA) Workshop* (IEEE, 2018) pp. 1–10.
 - [11] M. Baity-Jesi, L. Sagun, M. Geiger, S. Spigler, G. B. Arous, C. Cammarota, Y. LeCun, M. Wyart, and G. Biroli, in *International Conference on Machine Learning* (PMLR, 2018) pp. 314–323.
 - [12] M. Geiger, S. Spigler, S. d’Ascoli, L. Sagun, M. Baity-Jesi, G. Biroli, and M. Wyart, *Physical Review E* **100**, 12115 (2019).
 - [13] U. Simsekli, L. Sagun, and M. Gurbuzbalaban, in *International Conference on Machine Learning* (PMLR, 2019) pp. 5827–5837.
 - [14] Y. Feng and Y. Tu, *Proceedings of the National Academy of Sciences* **118** (2021), 10.1073/pnas.2015617118.
 - [15] K. He, X. Zhang, S. Ren, and J. Sun, in *Proceedings of the IEEE conference on computer vision and pattern recognition* (2016) pp. 770–778.
 - [16] P. Charbonneau, J. Kurchan, G. Parisi, P. Urbani, and F. Zamponi, *Nature Communications* **5**, 3725 (2014).
 - [17] H. J. Hwang, R. A. Riggelman, and J. C. Crocker, *Nature Materials* **15**, 1031 (2016).
 - [18] Y. Jin and H. Yoshino, *Nature Communications* **8**, 14935 (2017).
 - [19] P. Cao, M. P. Short, and S. Yip, *Proceedings of the National Academy of Sciences*, 201907317 (2019).
 - [20] I. Golding and E. C. Cox, *Physical Review Letters* **96**, 098102 (2006).
 - [21] I. Bronstein, Y. Israel, E. Kepten, S. Mai, Y. Shav-Tal, E. Barkai, and Y. Garini, *Physical Review Letters* **103**, 018102 (2009).
 - [22] V. Zaburdaev, S. Denisov, and J. Klafter, *Reviews of Modern Physics* **87**, 483 (2015).
 - [23] R. Metzler and J. Klafter, *Physics Reports* **339**, 1 (2000).
 - [24] T. H. Solomon, E. R. Weeks, and H. L. Swinney, *Physical Review Letters* **71**, 3975 (1993).
 - [25] G. M. Viswanathan, S. V. Buldyrev, S. Havlin, M. G. E. Da Luz, E. P. Raposo, and H. E. Stanley, *Nature* **401**, 911 (1999).
 - [26] P. Dieterich, R. Klages, R. Preuss, and A. Schwab, *Proceedings of the National Academy of Sciences* **105**, 459

- (2008).
- [27] L. G. Alves, D. B. Scariot, R. R. Guimarães, C. V. Nakamura, R. S. Mendes, and H. V. Ribeiro, *PLoS One* **11**, e0152092 (2016).
- [28] K. A. Sankararaman, S. De, Z. Xu, W. R. Huang, and T. Goldstein, in *Proceedings of the 37th International Conference on Machine Learning* (2020).
- [29] J. P. Nolan, *Univariate Stable Distributions: Models for Heavy Tailed Data* (Springer International Publishing, 2020) pp. 1–23.
- [30] J. Klafter and I. M. Sokolov, *First Steps in Random Walks: from Tools to Applications* (Oxford University Press, 2011).
- [31] M. F. Barnsley, R. L. Devaney, B. B. Mandelbrot, H.-O. Peitgen, D. Saupe, R. F. Voss, Y. Fisher, and M. McGuire, *The Science of Fractal Images* (Springer, 1988).
- [32] G. B. Sorkin, *Algorithmica* **6**, 367 (1991).
- [33] S. Hochreiter and J. Schmidhuber, *Neural Computation* **9**, 1 (1997).
- [34] L. Sagun, U. Evcı, V. U. Guney, Y. Dauphin, and L. Bottou, in *The International Conference on Learning Representations (ICLR) workshop* (2018).
- [35] C. Baldassi, F. Pittorino, and R. Zecchina, *Proceedings of the National Academy of Sciences* **117**, 161 (2020).
- [36] B. Ghorbani, S. Krishnan, and Y. Xiao, in *Proceedings of the 36th International Conference on Machine Learning*, *Proceedings of Machine Learning Research*, Vol. 97, edited by K. Chaudhuri and R. Salakhutdinov (PMLR, 2019) pp. 2232–2241.
- [37] P. Chaudhari, A. Choromanska, S. Soatto, Y. LeCun, C. Baldassi, C. Borgs, J. Chayes, L. Sagun, and R. Zecchina, *Journal of Statistical Mechanics: Theory and Experiment* **2019**, 124018 (2019).
- [38] J. Li, Q. Du, and C. Sun, *Pattern Recognition* **42**, 2460 (2009).
- [39] N. Douillet, *MATLAB Central File Exchange* (2020).
- [40] C. H. Martin and M. W. Mahoney, *arXiv:1810.01075* (2019).
- [41] J. Zhang, S. P. Karimireddy, A. Veit, S. Kim, S. J. Reddi, S. Kumar, and S. Sra, *arXiv:1912.03194* (2019).
- [42] A. Panigrahi, R. Somani, N. Goyal, and P. Netrapalli, *arXiv:1910.09626* (2019).
- [43] R. Metzler, J.-H. Jeon, A. G. Cherstvy, and E. Barkai, *Physical Chemistry Chemical Physics* **16**, 24128 (2014).
- [44] D. S. Grebenkov, *Physical Review E* **99**, 032133 (2019).
- [45] U. Şimşekli, O. Sener, G. Deligiannidis, and M. A. Erdogdu, “Hausdorff Dimension, Stochastic Differential Equations, and Generalization in Neural Networks,” (2020).
- [46] A. Wardak and P. Gong, *Physical Review Research* **3**, 013083 (2021).
- [47] J. A. Costa and A. O. Hero, in *European Signal Processing Conference* (2004) pp. 369–372.
- [48] E. Levina and P. J. Bickel, in *Neural Information Processing Systems* (2004) pp. 777–784.
- [49] S. Goldt, M. Mezard, F. Krzakala, and L. Zdeborova, *arXiv:1909.11500v3* (2019).
- [50] S. Santurkar, D. Tsipras, A. Ilyas, and A. Madry, in *Advances in Neural Information Processing Systems* (2018) pp. 2483–2493.